

# Knowledge Discovery and Data Mining

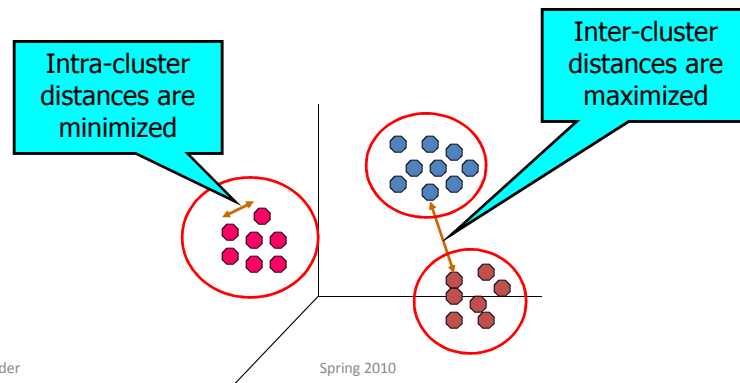
## Unit # 8

## Acknowledgement

- Most of the slides in this presentation are taken from course slides provided by
  - Han and Kimber (Data Mining Concepts and Techniques) and
  - Tan, Steinbach and Kumar (Introduction to Data Mining)
  - Several other online sources

## Cluster Analysis

- Finding groups of objects such that the objects in a group will be similar (or related) to one another and different from (or unrelated to) the objects in other groups



Sajjad Haider

Spring 2010

3

## Cluster Analysis (Cont'd)

- 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

Sajjad Haider

Spring 2010

4

## Google News

### iPhone activation headaches still trouble users

Computerworld - 1 hour ago  
 July 02, 2007 (Computerworld) -- It took Iain Gillott 47 hours to activate his iPhone after waiting in the Texas heat Friday afternoon to buy one. Most iPhone users thrilled but a few are iRate Reuters  
[Apple iPhone Arrives in the US](#) Techtree.com  
[Forbes](#) - [ZDNet](#) - [Ars Technica](#) - [Wired News](#)  
[all 562 news articles »](#)



- They didn't pick all 3,400,217 related articles by hand...
- Or Amazon.com
- Or Netflix...

### McCain Considers Ways to Reshape Campaign

Washington Post - 35 minutes ago  
 By Alec MacGillis Sen. John McCain's presidential campaign today announced widespread cutbacks and said it was considering whether to accept public campaign funds after another disappointing fundraising effort that has left the Arizona Republican with ...  
[McCain's Troubles Mount](#) New York Times  
[McCain Campaign Struggling, Reduces Staff](#) ABC News  
[CBS News](#) - [Reuters](#) - [Angus Reid Global Monitor](#) - [Sarasota Herald-Tribune](#)  
[all 291 news articles »](#)



## 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
- Earth-quake studies: Observed earth quake epicenters should be clustered along continent faults

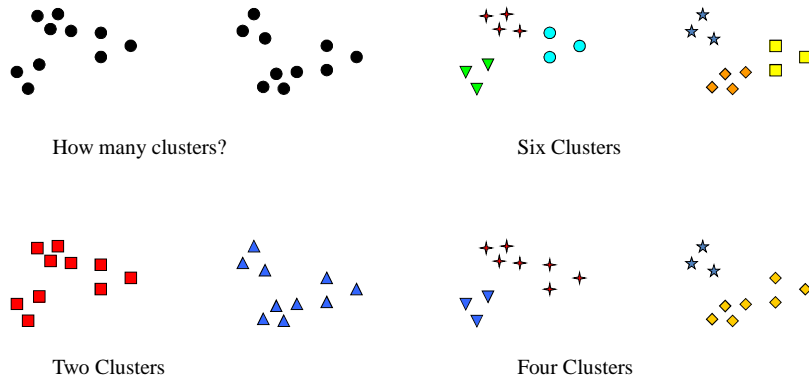
## Other less glamorous things...

- Hospital Records
- Scientific Imaging
  - Related genes, related stars, related sequences
- Market Research
  - Segmenting markets, product positioning
- Social Network Analysis
- Data mining
- Image segmentation...

## What is not Cluster Analysis?

- Supervised classification
  - Have class label information
- Simple segmentation
  - Dividing students into different registration groups alphabetically, by last name
- Results of a query
  - Groupings are a result of an external specification
- Graph partitioning
  - Some mutual relevance and synergy, but areas are not identical

## Notion of a Cluster can be Ambiguous



Sajjad Haider

Spring 2010

9

## Quality: What Is Good Clustering?

- A good clustering method will produce high quality clusters with
  - high intra-class similarity
  - low inter-class similarity
- The quality of a clustering result depends on both the similarity measure used by the method and its implementation
- The quality of a clustering method is also measured by its ability to discover some or all of the hidden patterns

Sajjad Haider

Spring 2010

10

## Types of Clusterings

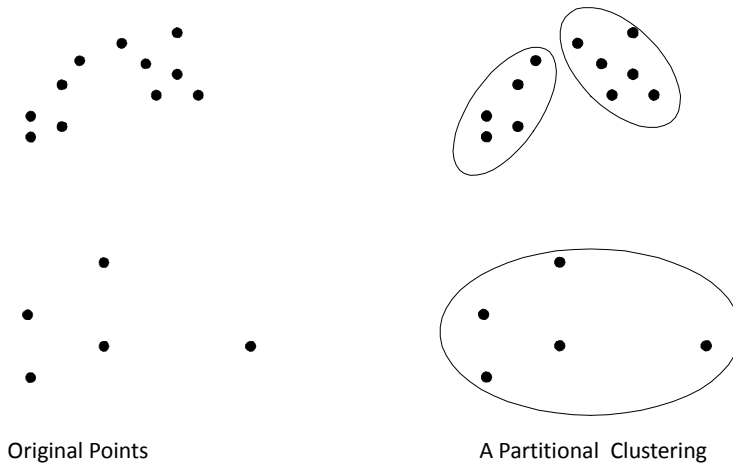
- A **clustering** is a set of clusters
- Important distinction between **hierarchical** and **partitional** sets of clusters
- **Partitional Clustering**
  - A division data objects into non-overlapping subsets (clusters) such that each data object is in exactly one subset
- **Hierarchical clustering**
  - A set of nested clusters organized as a hierarchical tree

Sajjad Haider

Spring 2010

11

## Partitional Clustering

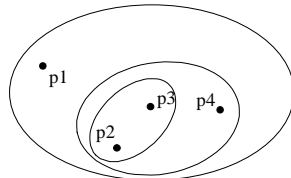


Sajjad Haider

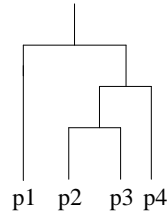
Spring 2010

12

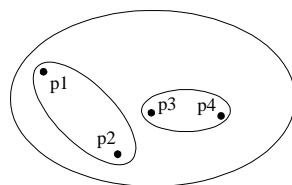
## Hierarchical Clustering



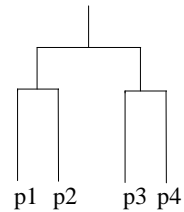
Traditional Hierarchical Clustering



Traditional Dendrogram



Non-traditional Hierarchical Clustering



Non-traditional Dendrogram

Sajjad Haider

Spring 2010

13

## 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.

Sajjad Haider

Spring 2010

14

## Typical Alternatives to Calculate the Distance between Clusters

- **Single link:** smallest distance between an element in one cluster and an element in the other, i.e.,  $\text{dis}(K_i, K_j) = \min(t_{ip}, t_{jq})$
- **Complete link:** largest distance between an element in one cluster and an element in the other, i.e.,  $\text{dis}(K_i, K_j) = \max(t_{ip}, t_{jq})$
- **Average:** avg distance between an element in one cluster and an element in the other, i.e.,  $\text{dis}(K_i, K_j) = \text{avg}(t_{ip}, t_{jq})$
- **Centroid:** distance between the centroids of two clusters, i.e.,  $\text{dis}(K_i, K_j) = \text{dis}(C_i, C_j)$
- **Medoid:** distance between the medoids of two clusters, i.e.,  $\text{dis}(K_i, K_j) = \text{dis}(M_i, M_j)$ 
  - Medoid: one chosen, centrally located object in the cluster

## What is Similarity?

The quality or state of being similar; likeness; resemblance; as, a similarity of features.  
**Webster's Dictionary**



Similarity is hard to define, but...  
*"We know it when we see it"*

The real meaning of similarity is a philosophical question. We will take a more pragmatic approach.



## Partitioning Algorithms: Basic Concept

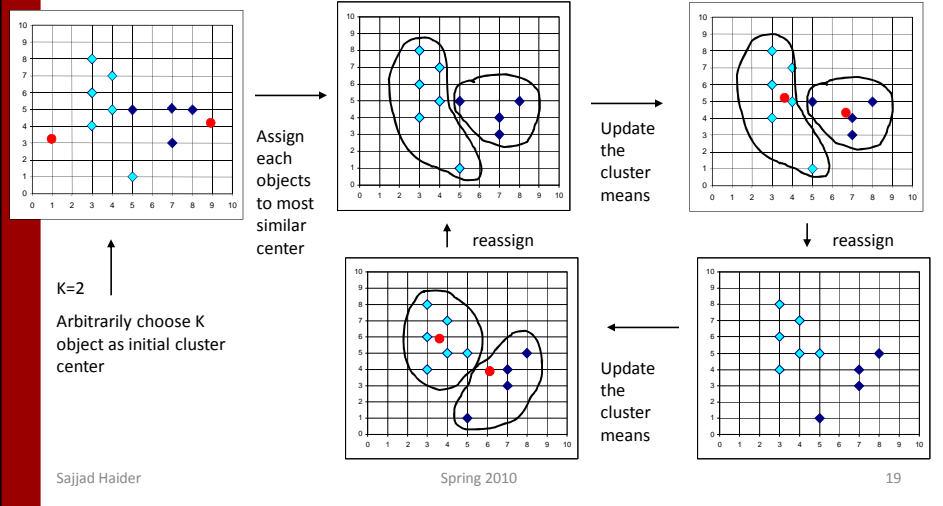
- Given a  $k$ , find a partition of  $k$  clusters that optimizes the chosen partitioning criterion
  - Global optimal: exhaustively enumerate all partitions
  - Heuristic methods: *k-means* and *k-medoids* algorithms
  - *k-means* (MacQueen'67): Each cluster is represented by the center of the cluster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

## 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

## The K-Means Clustering Method (Cont'd)

### • Example



## The K-Means Algorithm

1. Choose a value for  $K$ , the total number of clusters to be determined.
2. Choose  $K$  instances within the dataset at random. These are the initial cluster centers.
3. Use simple Euclidean distance to assign the remaining instances to their closest cluster center.
4. Use the instance in each cluster to calculate a new mean for each cluster.
5. If the new mean values are identical to the mean values of the previous iteration the process terminates. Otherwise, use the new means as cluster centers and repeat steps 3-5.

## Working of the K-Mean Algorithm

Instance #:	1	2	3	4	5	6
X:	1	1	2	2	3	3
Y:	1.5	4.5	1.5	3.5	2.5	6.0

- Let's pick Instances #1 and #3 as the initial centroids.

	Distance with Centroid1	Distance with Centroid2
• Instance #2	<b>3.00</b>	3.16
• Instance #4	2.24	<b>2.00</b>
• Instance #5	2.24	<b>1.41</b>
• Instance #6	6.02	<b>5.41</b>

- New centroids are (1, 3) and (2.5, 3.4)**

## Comments on the K-Means Method

- Strength:** *Relatively efficient:  $O(tkn)$ , where  $n$  is # objects,  $k$  is # clusters, and  $t$  is # iterations. Normally,  $k, t \ll n$ .*
- Comment:** Often terminates at a *local optimum*. The *global optimum* may be found using techniques such as: *deterministic annealing* and *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-convex shapes*

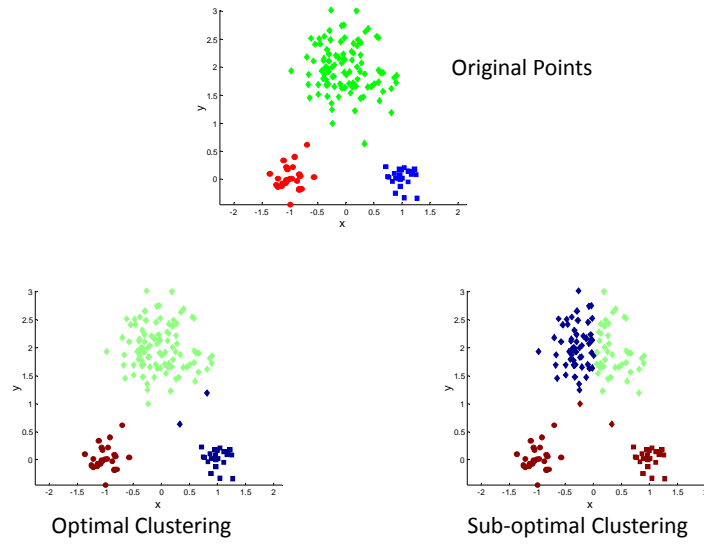
## Variations of the *K-Means* Method

- A few variants of the *k-means* which differ in
  - Selection of the initial *k* means
  - Dissimilarity calculations
  - Strategies to calculate cluster means
- 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

## K-Means Clustering – Details

- Initial centroids are often chosen randomly.
  - Clusters produced vary from one run to another.
- The centroid is (typically) the mean of the points in the cluster.
- 'Closeness' is measured by Euclidean distance, cosine similarity, correlation, etc.
- K-means will converge for common similarity measures mentioned above.
- Most of the convergence happens in the first few iterations.
  - Often the stopping condition is changed to 'Until relatively few points change clusters'

## Two different K-means Clusterings

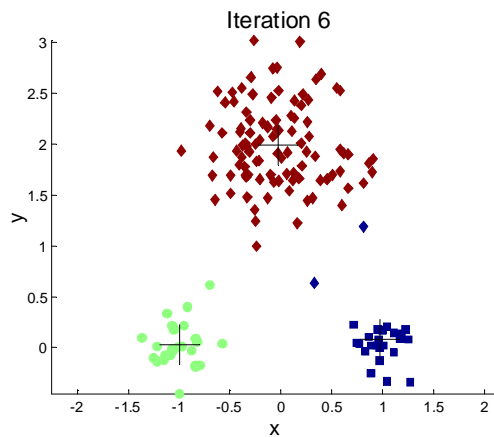


Sajjad Haider

Spring 2010

25

## Importance of Choosing Initial Centroids

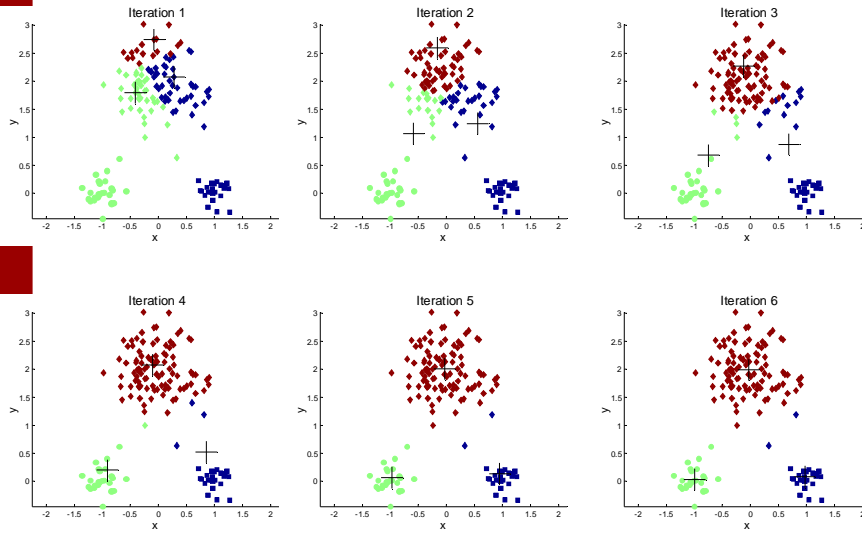


Sajjad Haider

Spring 2010

26

## Importance of Choosing Initial Centroids

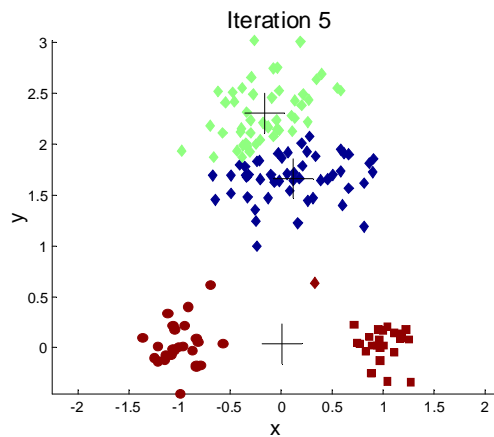


Sajjad Haider

Spring 2010

27

## Importance of Choosing Initial Centroids ...

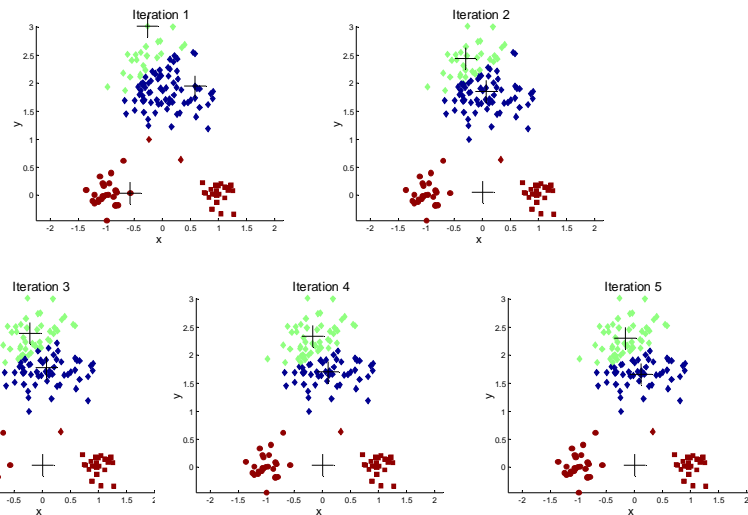


Sajjad Haider

Spring 2010

28

## Importance of Choosing Initial Centroids ...



Sajjad Haider

Spring 2010

29

## Evaluating K-means Clusters

- Most common measure is Sum of Squared Error (SSE)
  - For each point, the error is the distance to the nearest cluster
  - To get SSE, we square these errors and sum them.

$$SSE = \sum_{i=1}^K \sum_{x \in C_i} dist^2(m_i, x)$$

- $x$  is a data point in cluster  $C_i$  and  $m_i$  is the representative point for cluster  $C_i$ 
  - can show that  $m_i$  corresponds to the center (mean) of the cluster
- Given two clusters, we can choose the one with the smallest error
- One easy way to reduce SSE is to increase  $K$ , the number of clusters
  - A good clustering with smaller  $K$  can have a lower SSE than a poor clustering with higher  $K$

Sajjad Haider

Spring 2010

30

## Solutions to Initial Centroids Problem

- Multiple runs
  - Helps, but probability is not on your side
- Sample and use hierarchical clustering to determine initial centroids
- Select more than  $k$  initial centroids and then select among these initial centroids
  - Select most widely separated
- Postprocessing
- Bisecting K-means
  - Not as susceptible to initialization issues

Sajjad Haider

Spring 2010

31

## Pre-processing and Post-processing

- Pre-processing
  - Normalize the data
  - Eliminate outliers
- Post-processing
  - Eliminate small clusters that may represent outliers
  - Split 'loose' clusters, i.e., clusters with relatively high SSE
  - Merge clusters that are 'close' and that have relatively low SSE

Sajjad Haider

Spring 2010

32



## Bisecting K-means

- Bisecting K-means algorithm
  - Variant of K-means that can produce a partitional or a hierarchical clustering

- 
- 1: Initialize the list of clusters to contain the cluster containing all points.
  - 2: **repeat**
  - 3: Select a cluster from the list of clusters
  - 4: **for**  $i = 1$  to *number\_of\_iterations* **do**
  - 5: Bisect the selected cluster using basic K-means
  - 6: **end for**
  - 7: Add the two clusters from the bisection with the lowest SSE to the list of clusters.
  - 8: **until** Until the list of clusters contains  $K$  clusters
- 

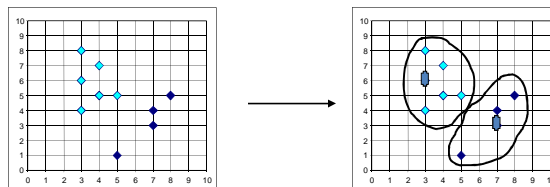
Sajjad Haider

Spring 2010

33

## 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.



Sajjad Haider

Spring 2010

34

## Limitations of K-means

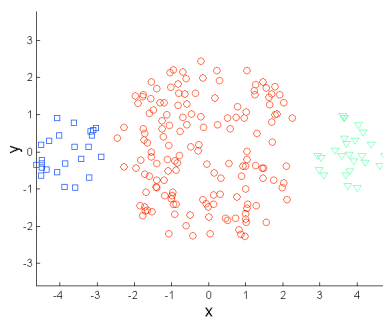
- K-means has problems when clusters are of differing
  - Sizes
  - Densities
  - Non-globular shapes
- K-means has problems when the data contains outliers.

Sajjad Haider

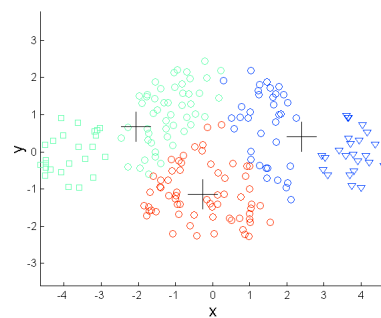
Spring 2010

35

## Limitations of K-means: Differing Sizes



Original Points



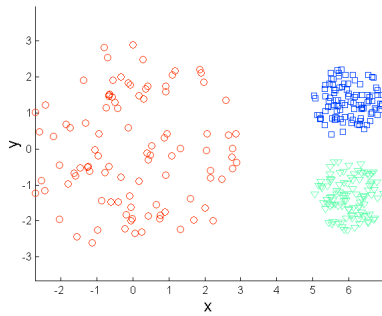
K-means (3 Clusters)

Sajjad Haider

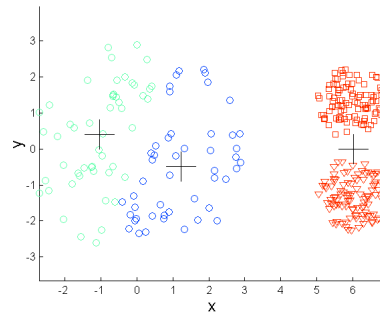
Spring 2010

36

## Limitations of K-means: Differing Density

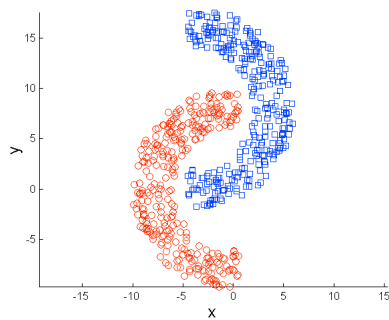


Original Points

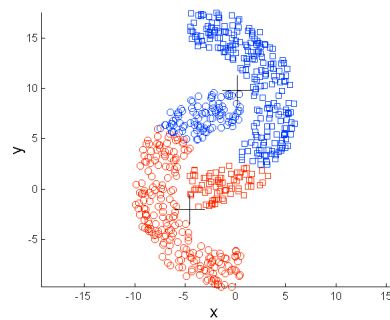


K-means (3 Clusters)

## Limitations of K-means: Non-globular Shapes

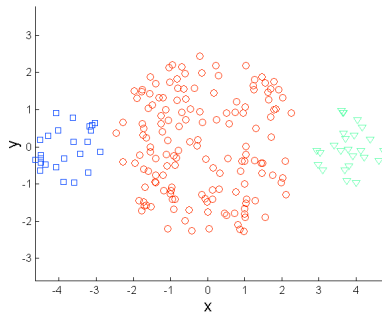


Original Points

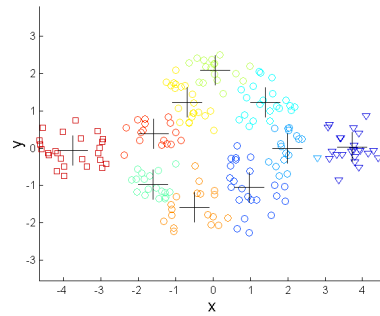


K-means (2 Clusters)

## Overcoming K-means Limitations



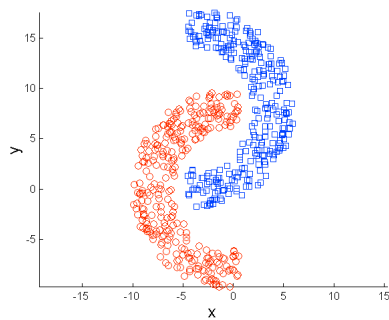
Original Points



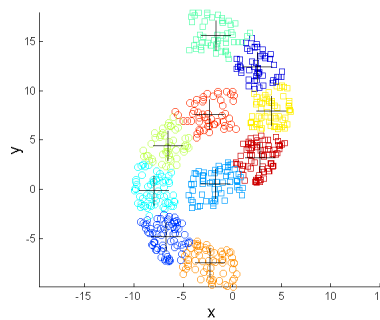
K-means Clusters

One solution is to use many clusters.  
Find parts of clusters, but need to put together.

## Overcoming K-means Limitations



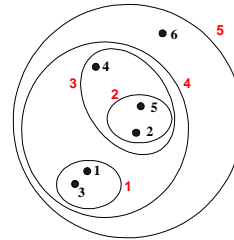
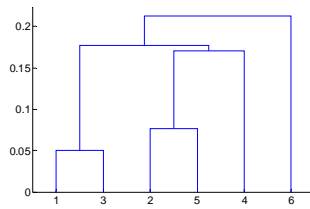
Original Points



K-means Clusters

## Hierarchical Clustering

- Produces a set of nested clusters organized as a hierarchical tree
- Can be visualized as a dendrogram
  - A tree like diagram that records the sequences of merges or splits



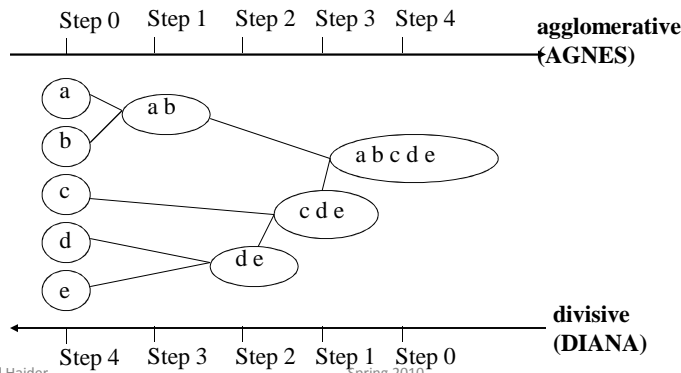
Sajjad Haider

Spring 2010

41

## Hierarchical Clustering (Cont'd)

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



Sajjad Haider

Spring 2010

42

## Strengths of Hierarchical Clustering

- Do not have to assume any particular number of clusters
  - Any desired number of clusters can be obtained by 'cutting' the dendrogram at the proper level
- They may correspond to meaningful taxonomies
  - Example in biological sciences (e.g., animal kingdom, phylogeny reconstruction, ...)

## Hierarchical Clustering

- Two main types of hierarchical clustering
  - Agglomerative:
    - Start with the points as individual clusters
    - At each step, merge the closest pair of clusters until only one cluster (or k clusters) left
  - Divisive:
    - Start with one, all-inclusive cluster
    - At each step, split a cluster until each cluster contains a point (or there are k clusters)
- Traditional hierarchical algorithms use a similarity or distance matrix
  - Merge or split one cluster at a time

## Agglomerative Clustering Algorithm

- More popular hierarchical clustering technique
- Basic algorithm is straightforward
  1. Compute the proximity matrix
  2. Let each data point be a cluster
  3. **Repeat**
  4. Merge the two closest clusters
  5. Update the proximity matrix
  6. **Until** only a single cluster remains
- Key operation is the computation of the proximity of two clusters
  - Different approaches to defining the distance between clusters distinguish the different algorithms

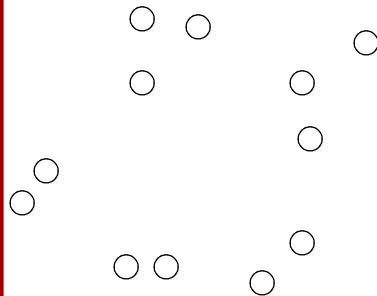
Sajjad Haider

Spring 2010

45

## Starting Situation

- Start with clusters of individual points and a proximity matrix



	p1	p2	p3	p4	p5	...
p1						
p2						
p3						
p4						
p5						
...						

Proximity Matrix



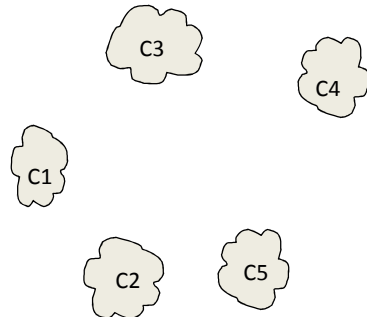
Sajjad Haider

Spring 2010

46

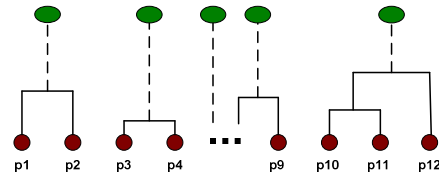
## Intermediate Situation

- After some merging steps, we have some clusters



	C1	C2	C3	C4	C5
C1					
C2					
C3					
C4					
C5					

Proximity Matrix



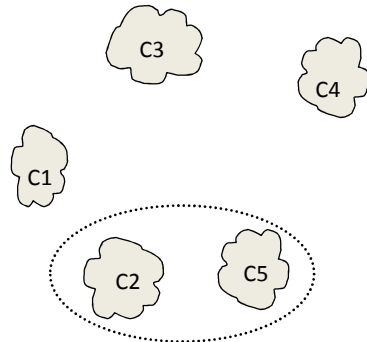
Sajjad Haider

Spring 2010

47

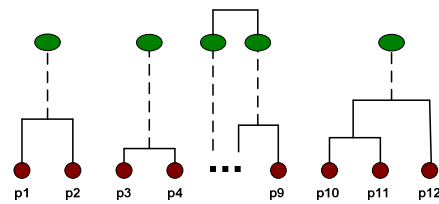
## Intermediate Situation

- We want to merge the two closest clusters (C2 and C5) and update the proximity matrix.



	C1	C2	C3	C4	C5
C1					
C2					
C3					
C4					
C5					

Proximity Matrix



Sajjad Haider

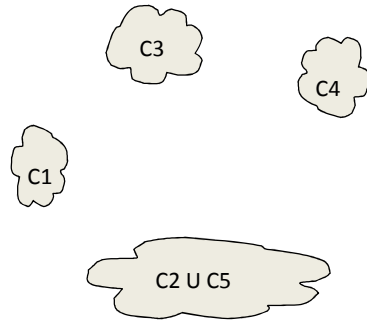
Spring 2010

48



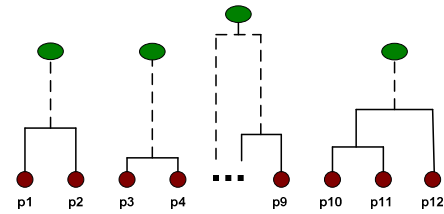
## After Merging

- The question is "How do we update the proximity matrix?"



		C1	C2 U C5	C3	C4
C1			?		
C2 U C5	?	?	?	?	?
C3			?		
C4			?		

Proximity Matrix

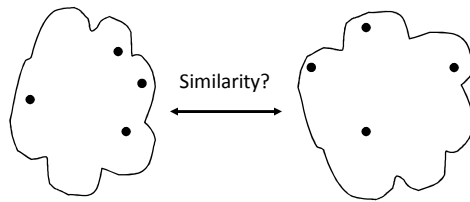


Sajjad Haider

Spring 2010

49

## How to Define Inter-Cluster Similarity



- MIN
- MAX
- Group Average
- Distance Between Centroids

	p1	p2	p3	p4	p5	...
p1						
p2						
p3						
p4						
p5						
.						
.						

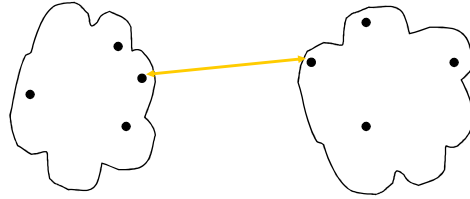
Proximity Matrix

Sajjad Haider

Spring 2010

50

## How to Define Inter-Cluster Similarity (Cont'd)

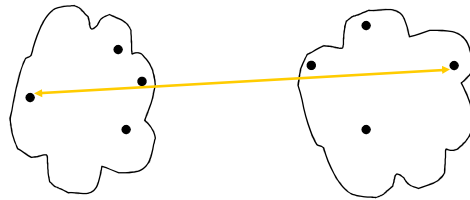


- MIN
- MAX
- Group Average
- Distance Between Centroids

	p1	p2	p3	p4	p5	...
p1						
p2						
p3						
p4						
p5						

Proximity Matrix

## How to Define Inter-Cluster Similarity (Cont'd)

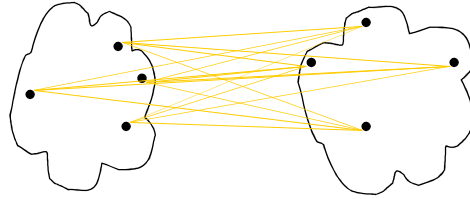


- MIN
- MAX
- Group Average
- Distance Between Centroids

	p1	p2	p3	p4	p5	...
p1						
p2						
p3						
p4						
p5						

Proximity Matrix

## How to Define Inter-Cluster Similarity (Cont'd)

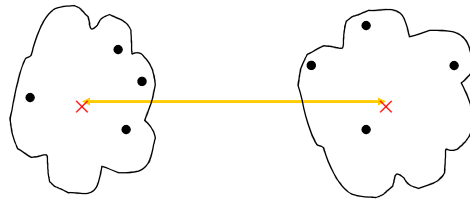


- MIN
- MAX
- Group Average
- Distance Between Centroids

	p1	p2	p3	p4	p5	...
p1						
p2						
p3						
p4						
p5						

Proximity Matrix

## How to Define Inter-Cluster Similarity (Cont'd)



- MIN
- MAX
- Group Average
- Distance Between Centroids

	p1	p2	p3	p4	p5	...
p1						
p2						
p3						
p4						
p5						

Proximity Matrix

## Single, Complete and Average Linkage

- In *single-linkage* clustering, the distance between one cluster and another cluster is equal to the shortest distance from any member of one cluster to any member of the other cluster:  $D(c_i, c_j) = \min d(a, b) \ a \in c_i, b \in c_j$ . It is obvious that:

$$D(c_k, c_l) = \min \{D(c_i, c_l), D(c_j, c_l)\} \text{ for } c_k = c_i \cup c_j$$

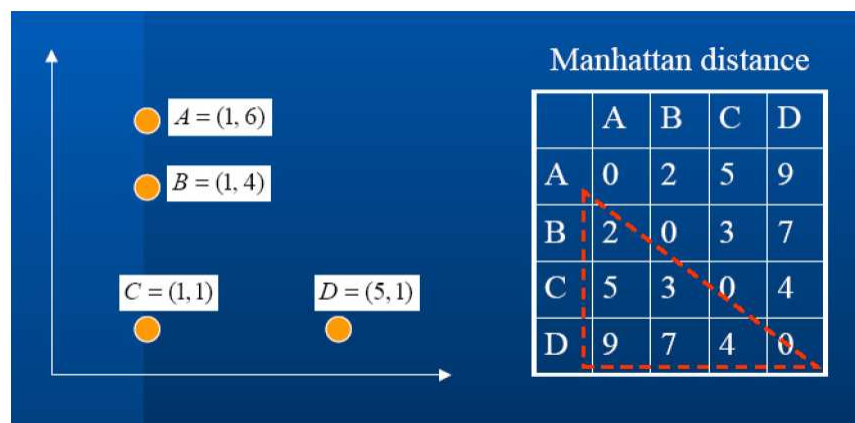
- In *complete-linkage* clustering, the distance between one cluster and another cluster is equal to the greatest distance from any member of one cluster to any member of the other cluster:  $D(c_i, c_j) = \max d(a, b) \ a \in c_i, b \in c_j$ .

- In *average-linkage* clustering, the distance between one cluster and another cluster is equal to the average distance from any member of one cluster to any member of the other cluster:

$$D(c_i, c_j) = \frac{1}{|c_i||c_j|} \sum_{a \in c_i, b \in c_j} d(a, b) \text{ . It is obvious that}$$

$$D(c_k, c_l) = \frac{|c_i|}{|c_k|} D(c_i, c_l) + \frac{|c_j|}{|c_k|} D(c_j, c_l) \text{ for } c_k = c_i \cup c_j$$

## Example: Distance Computation



## Example: Single Linkage

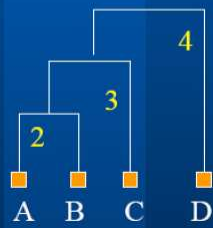
### Single linkage



$$\begin{aligned} \text{dist}((A, B), C) &= \min\{\text{dist}(A, C), \text{dist}(B, C)\} \\ &= \min\{5, 3\} = 3 \end{aligned}$$

$$\begin{aligned} \text{dist}((A, B), D) &= \min\{\text{dist}(A, D), \text{dist}(B, D)\} \\ &= \min\{9, 7\} = 7 \end{aligned}$$

$$\text{dist}(C, D) = 4$$

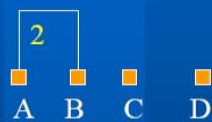


$$\begin{aligned} \text{dist}((A, B, C), D) &= \min\{\text{dist}((A, B), D), \text{dist}(C, D)\} \\ &= \min\{7, 4\} = 4 \end{aligned}$$

Spring 2010

## Example: Average Linkage

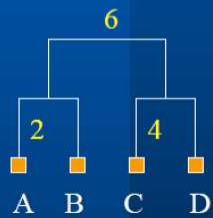
### Average linkage



$$\begin{aligned} \text{dist}((A, B), C) &= \text{avg}\{\text{dist}(A, C), \text{dist}(B, C)\} \\ &= (5+3)/2 = 4 \end{aligned}$$

$$\begin{aligned} \text{dist}((A, B), D) &= \text{avg}\{\text{dist}(A, D), \text{dist}(B, D)\} \\ &= (9+7)/2 = 8 \end{aligned}$$

$$\text{dist}(C, D) = 4$$



$$\begin{aligned} \text{dist}((C, D), (A, B)) &= \text{avg}\{\text{dist}(C, (A, B)), \text{dist}(D, (A, B))\} \\ &= (4+8)/2 = 6 \end{aligned}$$

Spring 2010

## Example: Complete Linkage

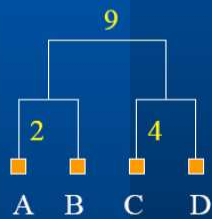
### Complete linkage



$$\begin{aligned} \text{dist}((A, B), C) &= \max \{ \text{dist}(A, C), \text{dist}(B, C) \} \\ &= \max \{ 5, 3 \} = 5 \end{aligned}$$

$$\begin{aligned} \text{dist}((A, B), D) &= \max \{ \text{dist}(A, D), \text{dist}(B, D) \} \\ &= \max \{ 9, 7 \} = 9 \end{aligned}$$

$$\text{dist}(C, D) = 4$$



$$\begin{aligned} \text{dist}((C, D), (A, B)) &= \max \{ \text{dist}(C, (A, B)), \text{dist}(D, (A, B)) \} \\ &= 9 \end{aligned}$$

Spring 2010

## Hierarchical Clustering: Problems and Limitations

- Once a decision is made to combine two clusters, it cannot be undone
- No objective function is directly minimized
- Different schemes have problems with one or more of the following:
  - Sensitivity to noise and outliers
  - Difficulty handling different sized clusters and convex shapes
  - Breaking large clusters

## Cluster Validity

- For supervised classification we have a variety of measures to evaluate how good our model is
  - Accuracy, precision, recall
- For cluster analysis, the analogous question is how to evaluate the “goodness” of the resulting clusters?
- But “clusters are in the eye of the beholder”!
- Then why do we want to evaluate them?
  - To avoid finding patterns in noise
  - To compare clustering algorithms
  - To compare two sets of clusters
  - To compare two clusters

## Measures of Cluster Validity

- Numerical measures that are applied to judge various aspects of cluster validity, are classified into the following three types.
  - **External Index:** Used to measure the extent to which cluster labels match externally supplied class labels.
    - Entropy
  - **Internal Index:** Used to measure the goodness of a clustering structure *without* respect to external information.
    - Sum of Squared Error (SSE)
  - **Relative Index:** Used to compare two different clusterings or clusters.
    - Often an external or internal index is used for this function, e.g., SSE or entropy

## Internal Measures: Cohesion and Separation

- **Cluster Cohesion:** Measures how closely related are objects in a cluster
  - Example: SSE
- **Cluster Separation:** Measure how distinct or well-separated a cluster is from other clusters
- Example: Squared Error
  - Cohesion is measured by the within cluster sum of squares (SSE)

$$WSS = \sum_i \sum_{x \in C_i} (x - m_i)^2$$

- Separation is measured by the between cluster sum of squares

$$BSS = \sum_i |C_i| (m - m_i)^2$$

- Where  $|C_i|$  is the size of cluster  $i$

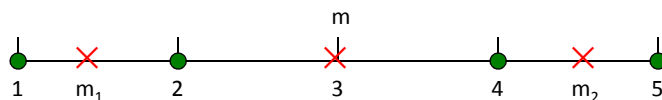
Sajjad Haider

Spring 2010

63

## Cohesion and Separation: Example

- Example: SSE
  - $BSS + WSS = \text{constant}$



K=1 cluster:

$$WSS = (1-3)^2 + (2-3)^2 + (4-3)^2 + (5-3)^2 = 10$$

$$BSS = 4 \times (3-3)^2 = 0$$

$$Total = 10 + 0 = 10$$

K=2 clusters:

$$WSS = (1-1.5)^2 + (2-1.5)^2 + (4-4.5)^2 + (5-4.5)^2 = 1$$

$$BSS = 2 \times (3-1.5)^2 + 2 \times (4.5-3)^2 = 9$$

$$Total = 1 + 9 = 10$$

Sajjad Haider

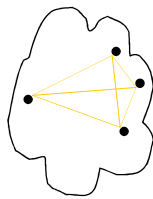
Spring 2010

64

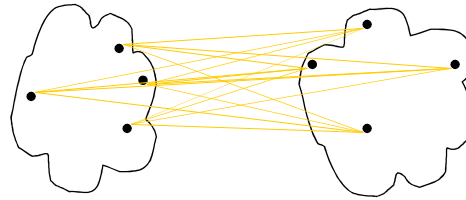


## Internal Measures: Cohesion and Separation (Cont'd)

- A proximity graph based approach can also be used for cohesion and separation.
  - Cluster cohesion is the sum of the weight of all links within a cluster.
  - Cluster separation is the sum of the weights between nodes in the cluster and nodes outside the cluster.



cohesion



separation

Sajjad Haider

Spring 2010

65

## External Measures of Cluster Validity: Entropy and Purity

Table 5.9. K-means Clustering Results for LA Document Data Set

Cluster	Entertainment	Financial	Foreign	Metro	National	Sports	Entropy	Purity
1	3	5	40	506	96	27	1.2270	0.7474
2	4	7	280	29	39	2	1.1472	0.7756
3	1	1	1	7	4	671	0.1813	0.9796
4	10	162	3	119	73	2	1.7487	0.4390
5	331	22	5	70	13	23	1.3976	0.7134
6	5	358	12	212	48	13	1.5523	0.5525
Total	354	555	341	943	273	738	1.1450	0.7203

**entropy** For each cluster, the class distribution of the data is calculated first, i.e., for cluster  $j$  we compute  $p_{ij}$ , the 'probability' that a member of cluster  $j$  belongs to class  $i$  as follows:  $p_{ij} = m_{ij}/m_j$ , where  $m_j$  is the number of values in cluster  $j$  and  $m_{ij}$  is the number of values of class  $i$  in cluster  $j$ . Then using this class distribution, the entropy of each cluster  $j$  is calculated using the standard formula  $e_j = \sum_{i=1}^L p_{ij} \log_2 p_{ij}$ , where the  $L$  is the number of classes. The total entropy for a set of clusters is calculated as the sum of the entropies of each cluster weighted by the size of each cluster, i.e.,  $e = \sum_{j=1}^K \frac{m_j}{m} e_j$ , where  $m_j$  is the size of cluster  $j$ ,  $K$  is the number of clusters, and  $m$  is the total number of data points.

**purity** Using the terminology derived for entropy, the purity of cluster  $j$ , is given by  $\text{purity}_j = \max_i p_{ij}$  and the overall purity of a clustering by  $\text{purity} = \sum_{j=1}^K \frac{m_j}{m} \text{purity}_j$ .

Sajjad Haider

Spring 2010

66

## Final Comment on Cluster Validity

“The validation of clustering structures is the most difficult and frustrating part of cluster analysis.

Without a strong effort in this direction, cluster analysis will remain a black art accessible only to those true believers who have experience and great courage.”

*Algorithms for Clustering Data, Jain and Dubes*

## Summary

- **Cluster analysis** groups objects based on their **similarity** and has wide applications
- Measure of similarity can be computed for **various types of data**
- Clustering algorithms can primarily be **categorized** into partitioning methods and hierarchical methods (there are many other less popular schemes too)
- **Outlier detection** and analysis are very useful for fraud detection, etc. and can be performed by statistical, distance-based or deviation-based approaches
- There are still lots of research issues on cluster analysis