

---

# Data Mining: Concepts and Techniques

## Cluster Analysis

Li Xiong

Slide credits: Jiawei Han and Micheline Kamber  
Tan, Steinbach, Kumar

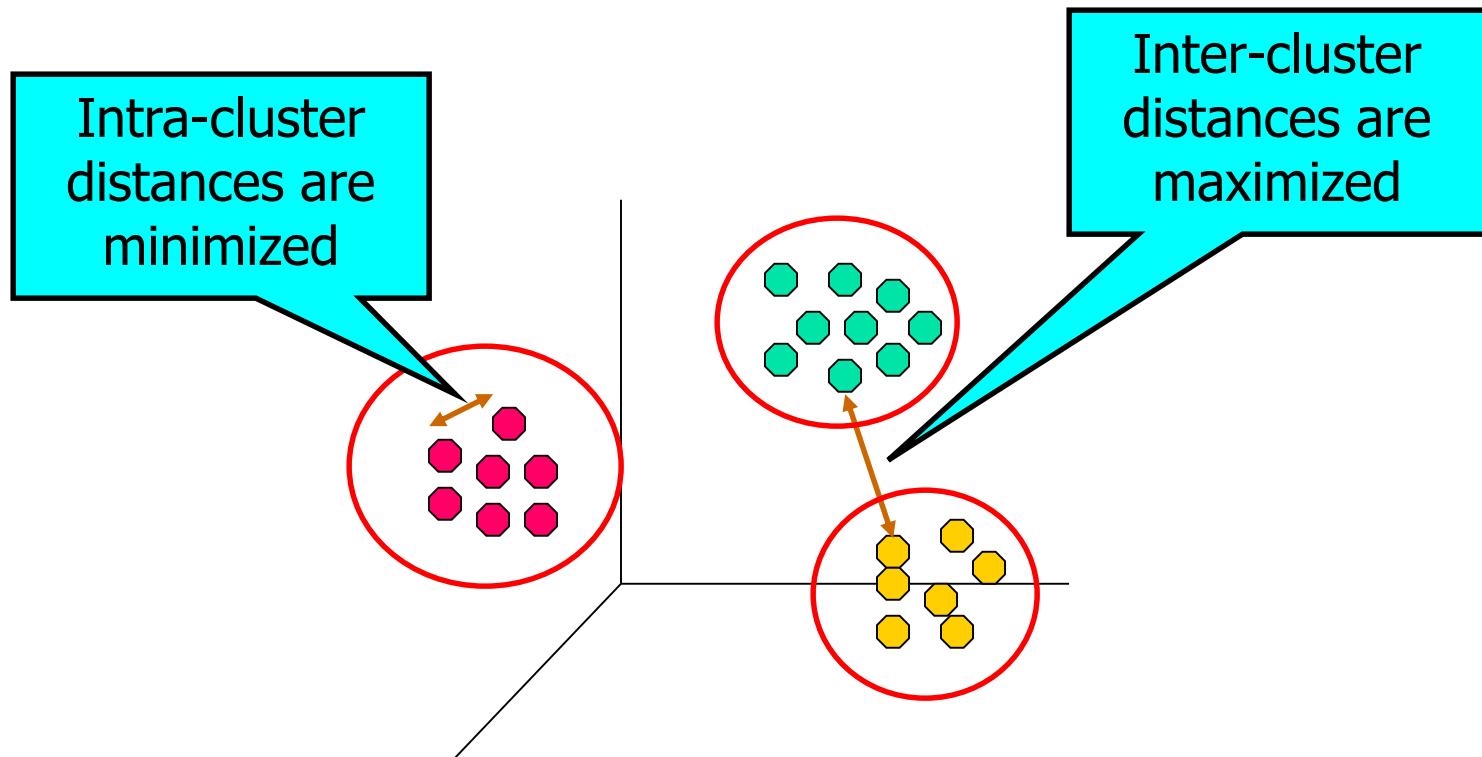
# Cluster Analysis

---

- Basic Concepts
- Similarity and distances
- Partitioning Methods
- Hierarchical Methods
- Density-Based Methods
- Probabilistic Methods
- Evaluation of Clustering

# What is Cluster Analysis?

- Finding groups of objects (clusters) – **given a notion of distance**
  - Objects **similar** to one another in the same group
  - Objects **different** from the objects in other groups
- Unsupervised learning: no predefined classes



# Machine Learning

- **Supervised:** Given input/output samples  $(X, y)$ , we learn a function  $f$  such that  $y = f(X)$ , which can be used on new data.
  - **Classification:**  $y$  is discrete (class labels).
  - **Regression:**  $y$  is continuous, e.g. linear regression.
- **Unsupervised:** Given only samples  $X$ , we compute a function  $f$  such that  $y = f(X)$  is “simpler”.
  - **Clustering:**  $y$  is discrete
  - **Dimension reduction:**  $y$  is continuous, e.g. matrix factorization

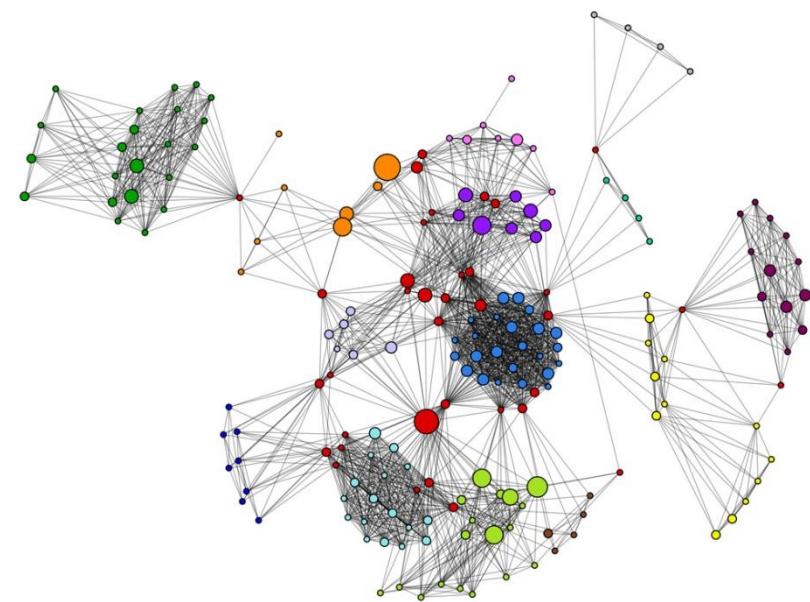
# Applications of Cluster Analysis

---

- As a stand-alone tool to get insight into data distribution
  - Cluster into groups – automatic classification
  - Finding k-nearest neighbors
  - Outlier detection
- As a preprocessing step for other algorithms
  - Data cleaning: missing data, noisy data
  - Data reduction
  - Data discretization

# Clustering Applications

- Marketing research
- Social network analysis



# Clustering Applications

## ■ WWW: Documents and search results clustering

The screenshot shows the Yippy search interface. At the top, there's a navigation bar with links for 'web', 'news', 'wikipedia', 'jobs', and 'more'. Below that is a search bar containing the query 'apple', with a 'Search' button and 'advanced preferences' link. The main area has tabs for 'clouds', 'sources', 'sites', and 'time', with 'clouds' currently selected. A sidebar on the left lists various search terms and their counts: iPad (42), Watch (38), OS (13), San Bernardino (54), Fight, FBI (49), Sales (29), Unlock, iPhone (38), Cook (29), App (25), Apple Mac (24), Pay (20), Apple And The Fbi (17), Valley (19), Reviews (19), Trump, Apple boycott (15), Market (18), Photos (17), Value (10), Repair (15), and Fruit (11). The main content area displays the top 505 results for the query 'apple'. Each result includes a title, a snippet of text, and links for sharing or viewing details. The results cover various topics like Apple products, CEO Tim Cook, and legal cases.

Top 505 results of at least 114,000,000 retrieved for the query [apple \(definition\)](#) ([details](#))

[Apple](#)

Apple leads the world in innovation with iPhone, iPad, Mac, Apple Watch, iOS, OS X, watchOS and more. Visit the site to learn, buy, and get support.  
<https://www.apple.com> - [cache] - Yippy Index IV

[iPad - Apple](#)

Discover the world of iPad. Introducing iPad Pro and the iPad mini 4. Visit the [Apple](#) site to learn, buy, and get support.  
[www.apple.com/ipad](http://www.apple.com/ipad) - [cache] - Yippy Index IV

[Apple's CEO: Complying with FBI demand 'bad for America'](#)

[www.boston.com/.../AUZDbduSBYe5lhMzvV4oqL/story.html](http://www.boston.com/.../AUZDbduSBYe5lhMzvV4oqL/story.html) - [cache] - Yippy News, Yippy News Archives

[Mark Hulbert: Apple's drop shows price of popularity](#)

[www.marketwatch.com/.../0\(MarketWatch.com - Newsletters & Research\)](http://www.marketwatch.com/.../0(MarketWatch.com - Newsletters & Research)) - [cache] - Yippy News, Yippy News Archives

[US cannot make Apple provide iPhone data in drug case, NY judge says | Fox News](#)

16 hours ago - The U.S. Justice Department cannot force [Apple](#) to provide the FBI with access to a locked iPhone's data in a Brooklyn drug case, a federal judge in New York ruled Monday.  
[www.foxnews.com/.../ake-apple-provide-iphone-data-in-drug-case-ny-judge-says.html](http://www.foxnews.com/.../ake-apple-provide-iphone-data-in-drug-case-ny-judge-says.html) - [cache] - Yippy News, Yippy News Archives

[Showdown On Gender Pay Equity In Silicon Valley: Shareholders Press Seven Tech Giants To Follow Lead Of Intel, Apple On Fair Treatment Of Women](#)

[www.prnewswire.com/.../d-of-intel-apple-on-fair-treatment-of-women-300229578.html](http://www.prnewswire.com/.../d-of-intel-apple-on-fair-treatment-of-women-300229578.html) - [cache] - Yippy News, Yippy News Archives

[Apple: iPhone Hacking Tool the 'Software Equivalent of Cancer'](#)

12 hours ago - [Apple](#) CEO Tim Cook doubled down on his stance against creating a backdoor software to unlock the iPhone of one of the San Bernadino terrorists, saying such a tool would be the software equivalent of cancer.  
[www.cnn.com/.../Apple-iPhone-Hacking-Tool-the-Software-Equivalent-of-Cancer](http://www.cnn.com/.../Apple-iPhone-Hacking-Tool-the-Software-Equivalent-of-Cancer) - [cache] - Yippy News, Yippy News Archives

[Apple backed by judge in new iPhone access fight - BBC News](#)

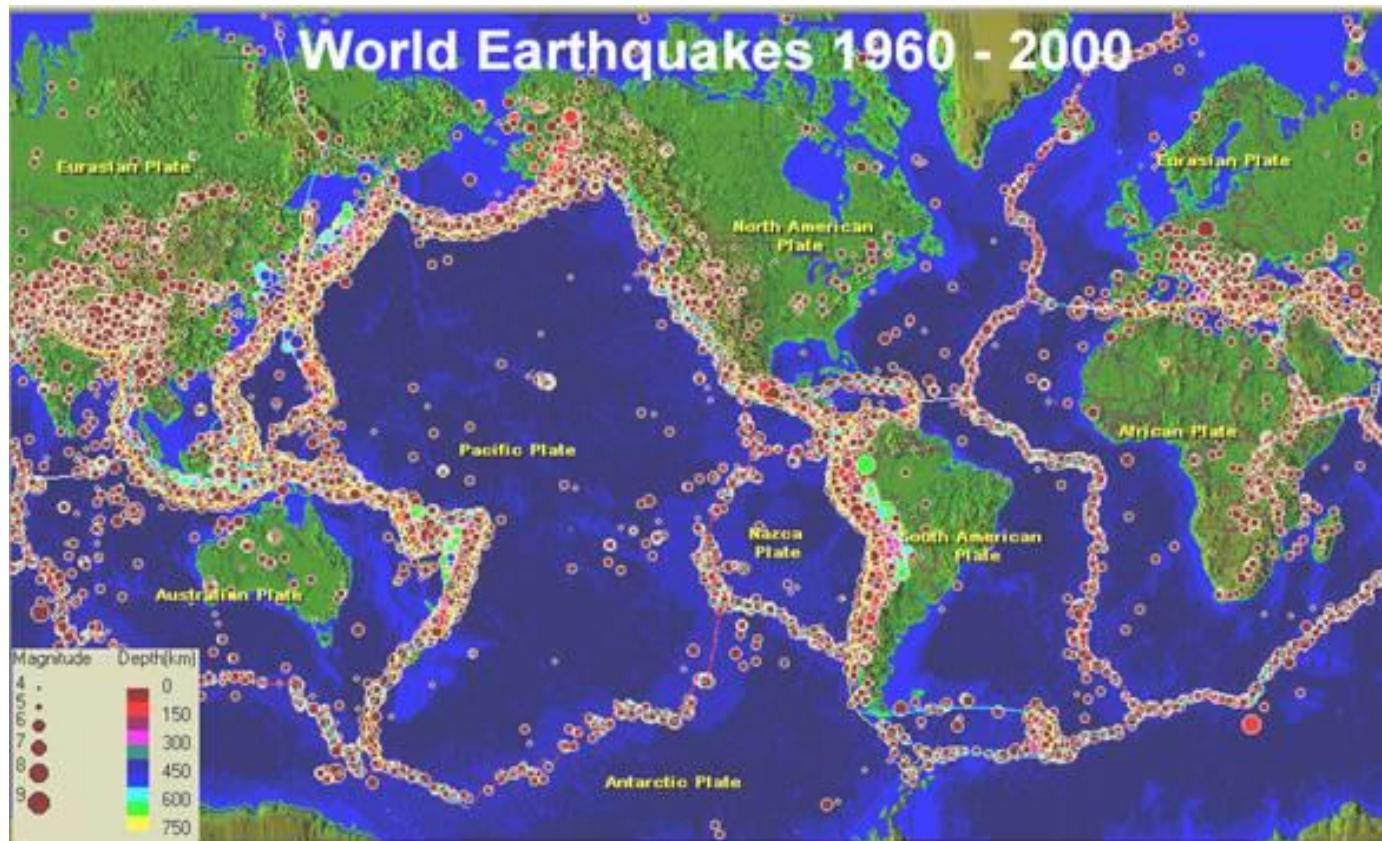
Mar 1, 2016 -  
[www.bbc.co.uk/news/world-us-canada-35692931](http://www.bbc.co.uk/news/world-us-canada-35692931) - [cache] - Yippy News, Yippy News

[The 20 bestselling mobile phones of all time](#)

[www.telegraph.co.uk/.../2016/01/26/the-20-bestselling-mobile-phones-of-all-time](http://www.telegraph.co.uk/.../2016/01/26/the-20-bestselling-mobile-phones-of-all-time) - [cache] - Yippy News, Yippy News Archives

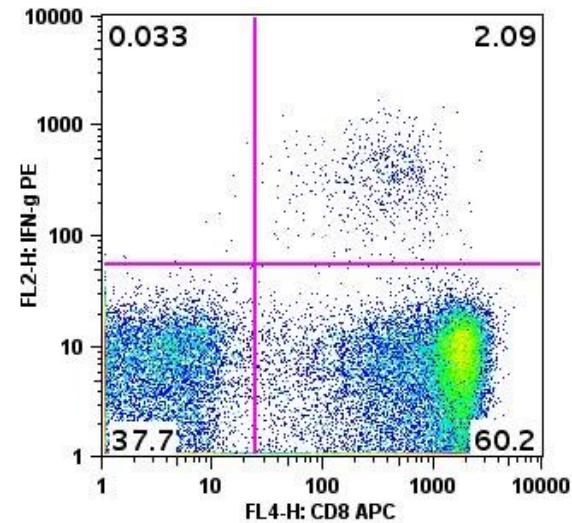
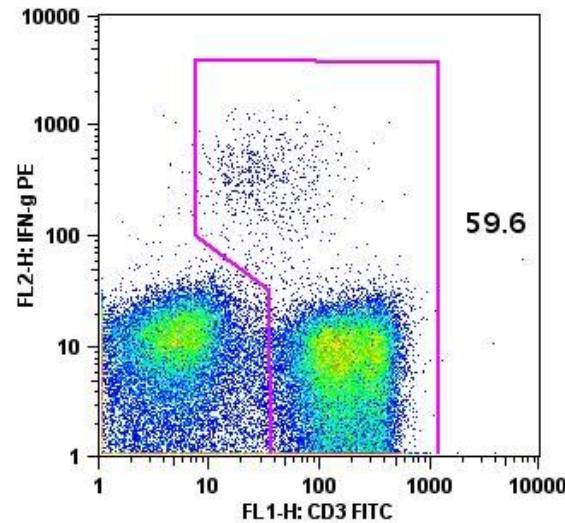
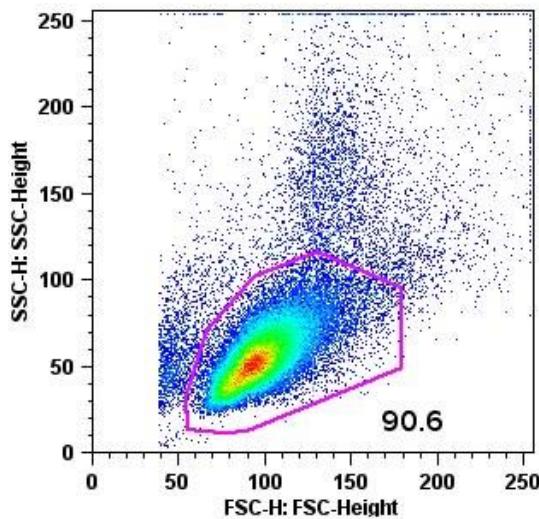
# Clustering Applications

- Earthquake studies



# Clustering Applications

- Bioinformatics: microarray data, flow cytometry data analysis, ...



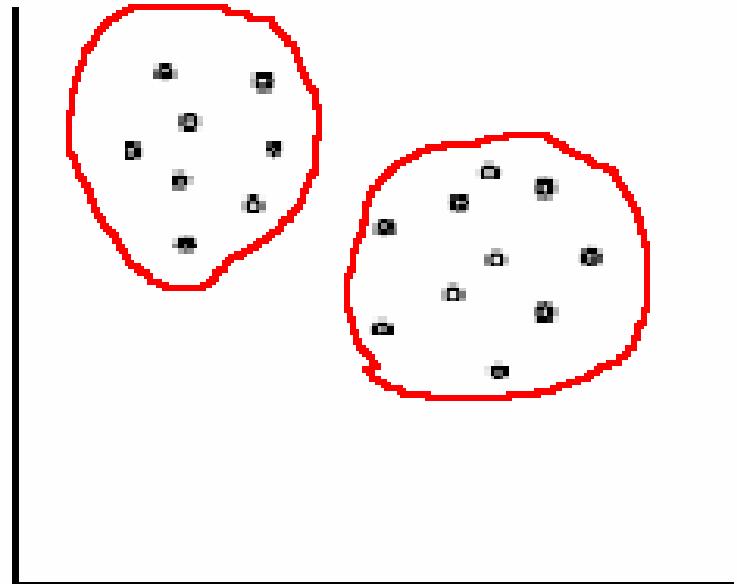
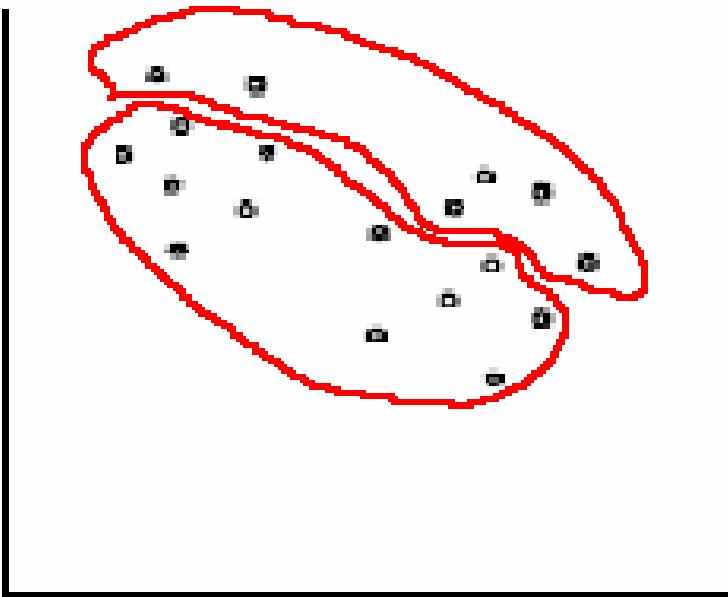
# Challenges of Clustering

---

- Quality
  - Noise and outliers
  - High dimensionality
- Scalability
  - High dimensionality
  - Large data
- Usability
  - Minimal input parameters
  - User-specified constraints

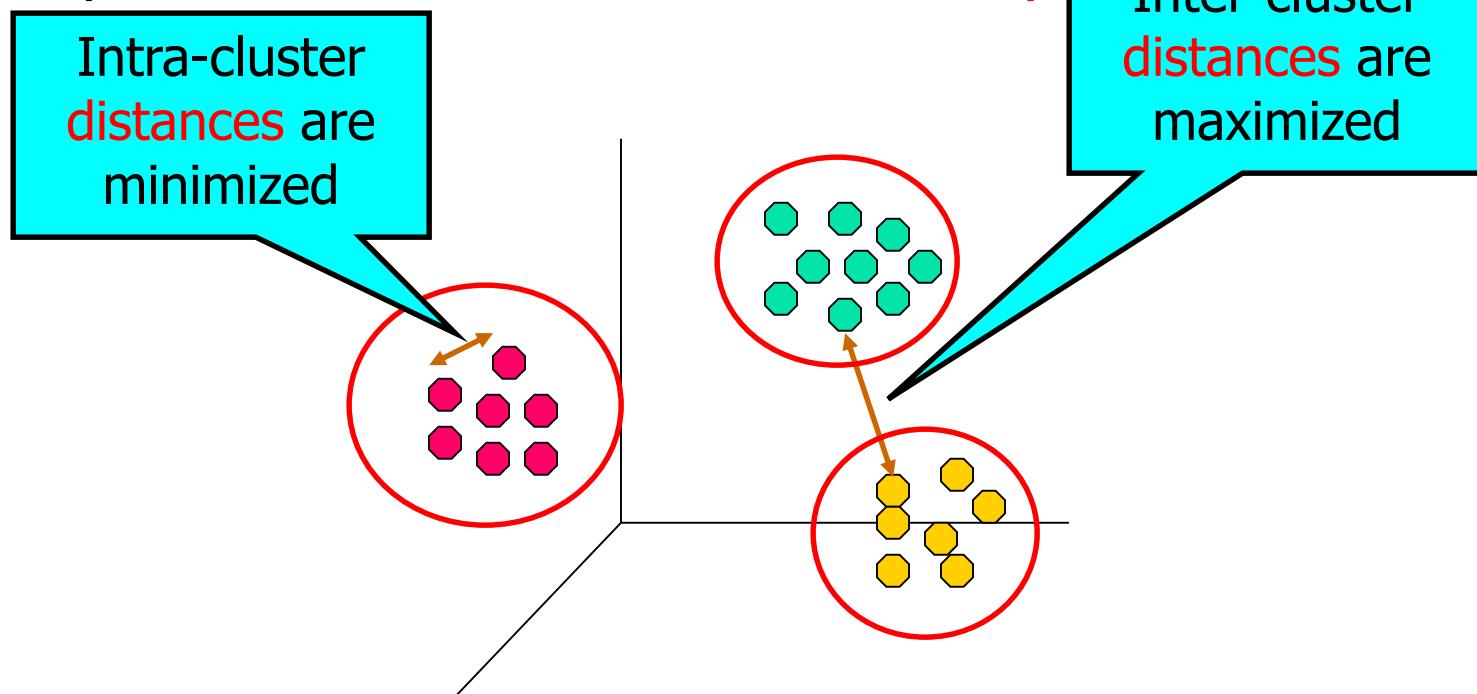
# Quality: What Is Good Clustering?

---



# Quality: What Is Good Clustering?

- Agreement with “ground truth”
- A good clustering will produce high quality clusters with
  - Homogeneity - high intra-class similarity
  - Separation - low inter-class similarity



# Cluster Analysis: Basic Concepts and Methods

---

- Cluster Analysis: Basic Concepts
- **Similarity and distances**
- Partitioning Methods
- Hierarchical Methods
- Density-Based Methods
- Probabilistic Methods
- Evaluation of Clustering

# Similarity/distance between data objects

---

## Data objects

- **as points:** distance between points
- **as vectors:** cosine between vectors
- **as random variables:** correlation
- **as sets:** Jaccard distance between sets
- **as strings:** Hamming distance

# Distance between two data points

- Euclidean distance

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

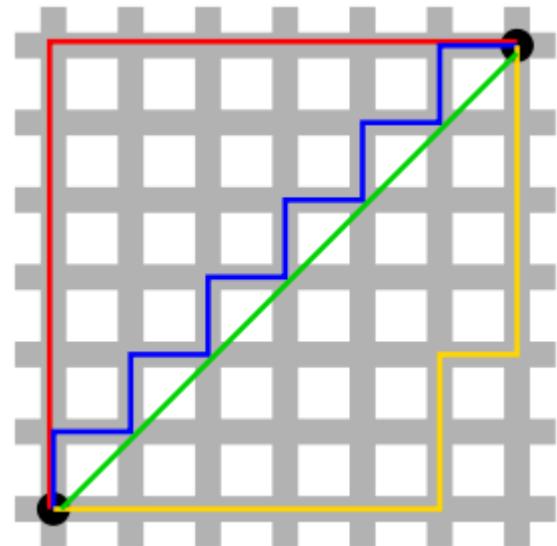
- Manhattan distance

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

- Minkowski distance

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

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



# Distance between two attributes values

---

- To compute  $|x_{if} - x_{jf}|$ 
  - $f$  is numeric (interval or ratio scale)
    - Scaling issues -> normalization
  - $f$  is ordinal
    - Mapping by rank 
$$z_{if} = \frac{r_{if} - 1}{M_f - 1}$$
  - $f$  is nominal
    - Mapping function
$$|x_{if} - x_{jf}| = 0 \text{ if } x_{if} = x_{jf}, \text{ or } 1 \text{ otherwise}$$
    - Hamming distance (edit distance) for strings

# Normalization of attributes

---

- scaled attributes to fall within a small, specified range
- Min-max normalization:  $[min_A, max_A]$  to  $[new\_min_A, new\_max_A]$

$$v' = \frac{v - min_A}{max_A - min_A} (new\_max_A - new\_min_A) + new\_min_A$$

- Ex. Let income  $[\$12,000, \$98,000]$  normalized to  $[0.0, 1.0]$ . Then  
 $\$73,600$  is mapped to  $\frac{73,600 - 12,000}{98,000 - 12,000} (1.0 - 0) + 0 = 0.716$
- Z-score normalization ( $\mu$ : mean,  $\sigma$ : standard deviation):

$$v' = \frac{v - \mu_A}{\sigma_A}$$

- Ex. Let  $\mu = 54,000$ ,  $\sigma = 16,000$ . Then  $\frac{73,600 - 54,000}{16,000} = 1.225$

# Euclidean distance

- Euclidean distance may not be meaningful (counter intuitive) for high dimensional data, e.g. user movie ratings

3 3 3 3 3 3 3 3 3 3 3

1 1 1 1 1 1 1 1 1 1 1

vs

0 3 0 3 0 3 0 3 0 3 0 3

1 1 1 1 1 1 1 1 1 1 1

# Similarity/distance between data objects

---

## Data objects

- as points: distance between points
- as vectors: cosine between vectors
- as random variables: correlation
- as sets: Jaccard distance between sets
- as strings: Hamming distance

# Cosine similarity between two vectors

---

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

- Cosine measure  $\frac{X_i \bullet X_j}{\|X_i\| \cdot \|X_j\|}$
- From -1 to 1

# Cosine similarity

- Cosine similarity
  - Invariant to multiplicative scaling
  - Variant to additive scaling

3 3 3 3 3 3 6 6 6 6 6 6

1 1 1 1 1 1 4 4 4 4 4 4

vs

2 2 2 2 2 2 8 8 8 8 8 8

1 1 1 1 1 1 4 4 4 4 4 4

# Correlation between two random variables (numerical data)

---

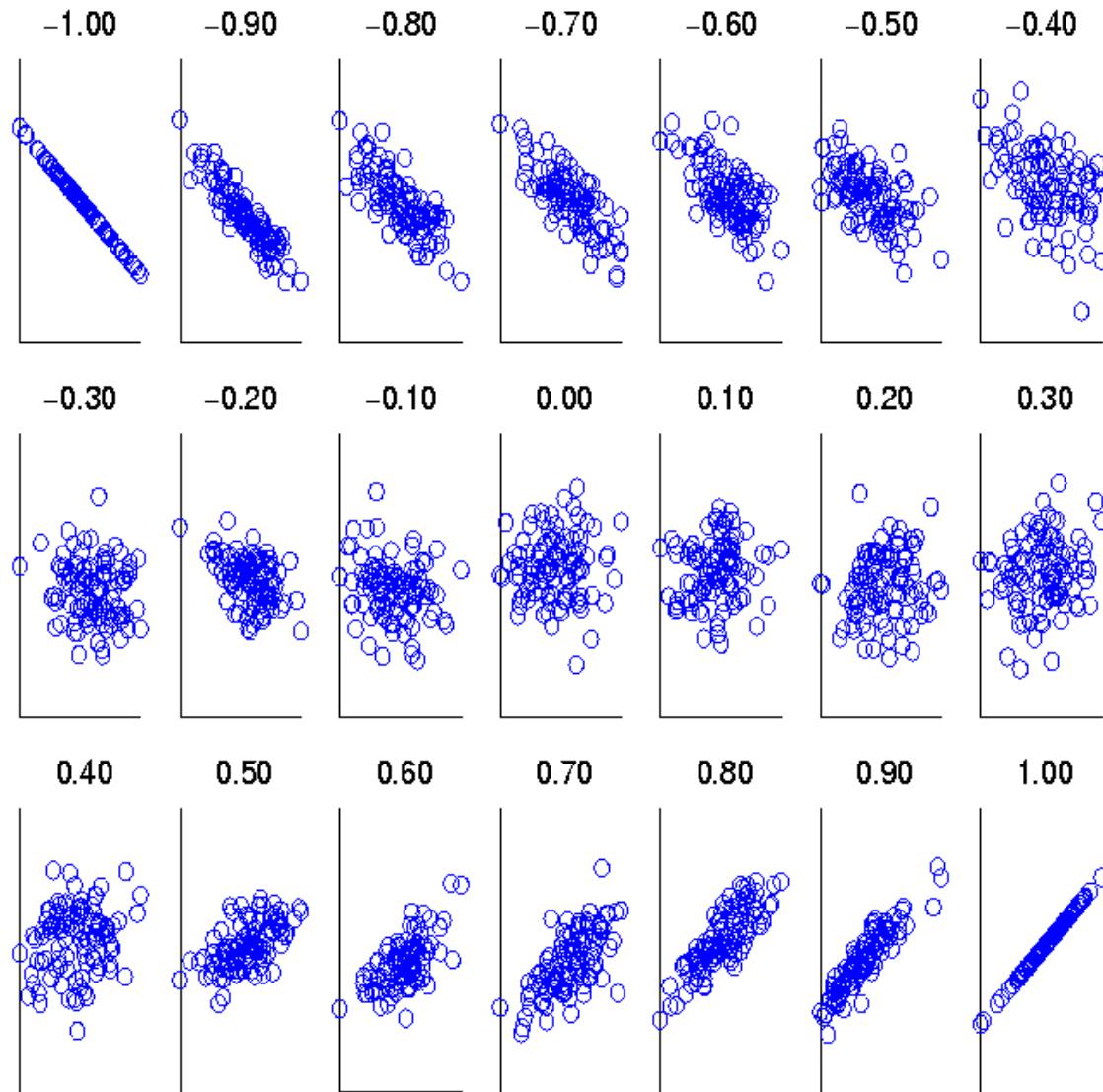
- Correlation coefficient (also called **Pearson's product moment coefficient**)

$$r_{A,B} = \frac{\sum (A - \bar{A})(B - \bar{B})}{(n-1)\sigma_A\sigma_B} = \frac{\sum (AB) - n\bar{A}\bar{B}}{(n-1)\sigma_A\sigma_B}$$

where  $n$  is the number of tuples,  $\bar{A}$  and  $\bar{B}$  are the respective means of A and B,  $\sigma_A$  and  $\sigma_B$  are the respective standard deviation of A and B, and  $\Sigma(AB)$  is the sum of the AB dot-product.

- $r_{A,B} > 0$ , A and B are positively correlated (A's values increase as B's)
- $r_{A,B} = 0$ : independent
- $r_{A,B} < 0$ : negatively correlated

# Visualization of Correlation



Scatter plots showing  
the Pearson correlation  
from  $-1$  to  $1$ .

# Data object at a set

- For transaction data, document data
  - Shared items are more important to consider

0 1 1 1 1 1 1 1 1 1 1

vs

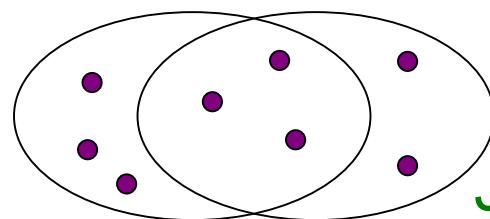
1 1 1 1 1 1 1 1 1 1 0

1 0 0 0 0 0 0 0 0 0 0

0 0 0 0 0 0 0 0 0 0 1

# Jaccard distance between two sets

- The **Jaccard similarity** of two **sets** is the size of their intersection divided by the size of their union:  
 $sim(C_1, C_2) = |C_1 \cap C_2| / |C_1 \cup C_2|$
- **Jaccard distance:**  $d(C_1, C_2) = 1 - |C_1 \cap C_2| / |C_1 \cup C_2|$



3 in intersection

8 in union

Jaccard similarity= 3/8

Jaccard distance = 5/8

# Cluster Analysis: Basic Concepts and Methods

---

- Cluster Analysis: Basic Concepts
- Similarity and distances
- **Partitioning Methods**
- Hierarchical Methods
- Density-Based Methods
- Probabilistic Methods
- Evaluation of Clustering

# Clustering Approaches

---

- Partitioning approach:
  - Construct various partitions and then evaluate them by some “goodness” criterion
  - Typical methods: k-means, k-medoids
- Hierarchical approach:
  - Create a hierarchical decomposition of the objects
  - Typical methods: Diana, Agnes
- Density-based approach:
  - Based on connectivity and density functions
  - Typical methods: DBSCAN
- Others

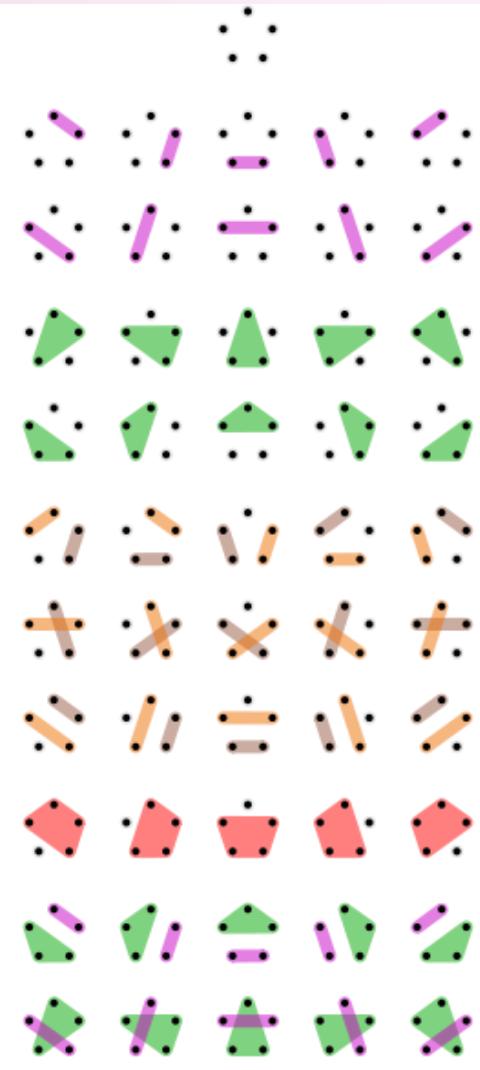
# Partitioning Algorithms: Basic Concept

---

- Partitioning method: Construct a partition of  $n$  objects (into  $k$  clusters), s.t. intracluster similarity maximized and intercluster similarity minimized
  - One objective: minimize the sum of squared distance from cluster centroid
- $$\sum_{i=1}^k \sum_{p \in C_i} (p - m_i)^2$$
- How to find optimal partition?

# Number of partitionings

---



# Number of partitionings

- Stirling partition number – number of ways to partition  $n$  objects into  $k$  non-empty subsets

$$\left\{ \begin{matrix} n+1 \\ k \end{matrix} \right\} = k \left\{ \begin{matrix} n \\ k \end{matrix} \right\} + \left\{ \begin{matrix} n \\ k-1 \end{matrix} \right\}$$

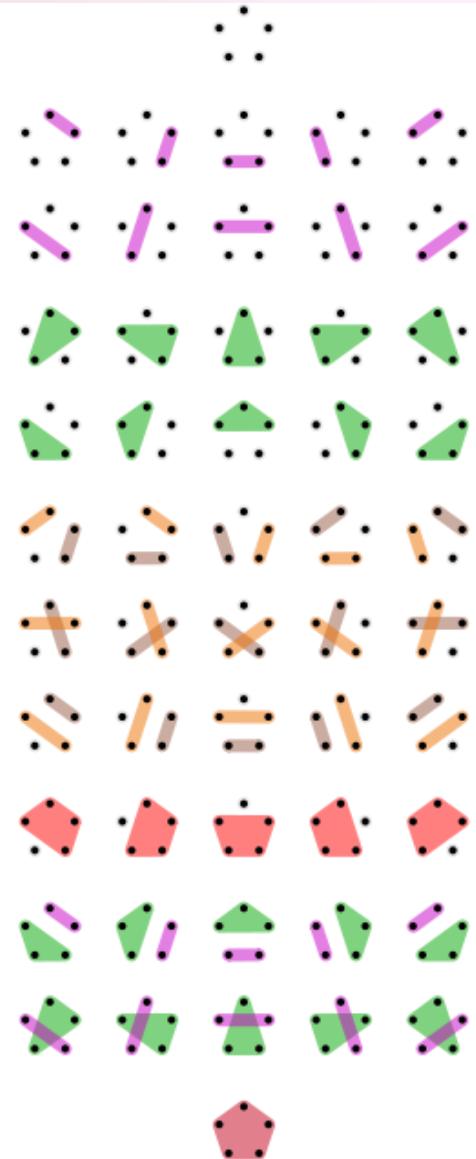
( $n=5, k=1, 2, 3, 4, 5$ ): 1, 15, 25, 10, 1

( $n=10, k=1, 2, 3, 4, 5, \dots$ ): 1, 511, 9330, 34105, 42525, ...

- Bell numbers – number of ways to partition  $n$  objects

$$B_n = \sum_{k=0}^n \left\{ \begin{matrix} n \\ k \end{matrix} \right\}.$$

( $n = 0, 1, 2, 3, 4, 5, \dots$ ): 1, 1, 2, 5, 15, 52, 203, 877, 4140, 21147, 115975, 678570, 4213597, 27644437, 190899322, 1382958545, 10480142147, 82864869804, 682076806159, 5832742205057, ...



# Partitioning Algorithms: Basic Concept

---

- Partitioning method: Construct a partition of  $n$  objects into  $k$  clusters, s.t. intracluster similarity maximized and intercluster similarity minimized
  - One objective: minimize the sum of squared distance from cluster centroid
- $$\sum_{i=1}^k \sum_{p \in C_i} (p - m_i)^2$$
- Heuristic methods: *k-means* and *k-medoids* algorithms
    - *k-means* (Lloyd'57, 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

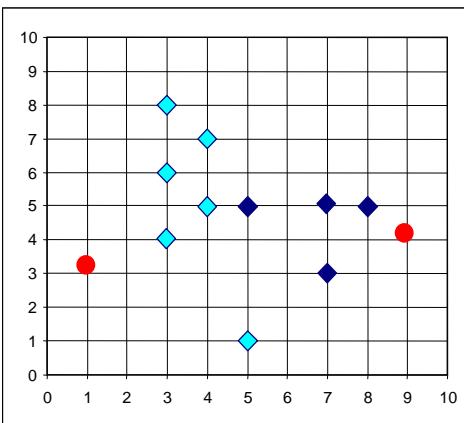
# *K-Means* Clustering: Lloyd Algorithm

---

- Given  $k$ , and randomly choose  $k$  initial cluster centers
- Partition objects into  $k$  nonempty subsets by assigning each object to the cluster with the **nearest** centroid
- Update centroid, i.e. ***mean point*** of the cluster
- Go back to Step 2, stop when no more new assignment and centroids do not change

# The *K*-Means Clustering Method

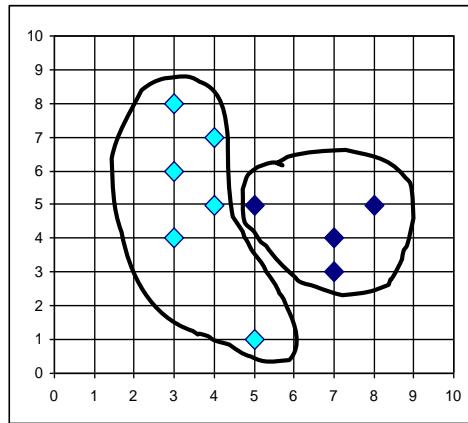
## ■ Example



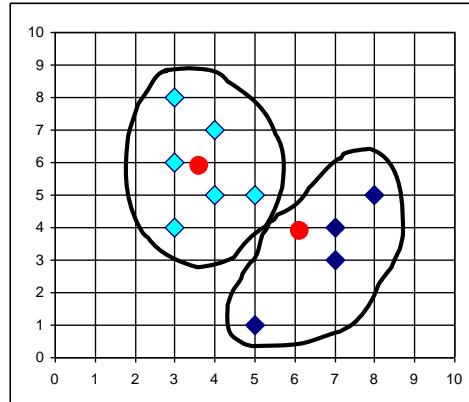
K=2

Arbitrarily choose K object as initial cluster center

Assign each objects to most similar center

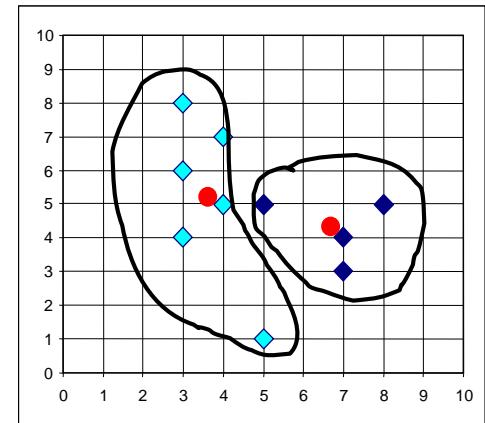


reassign

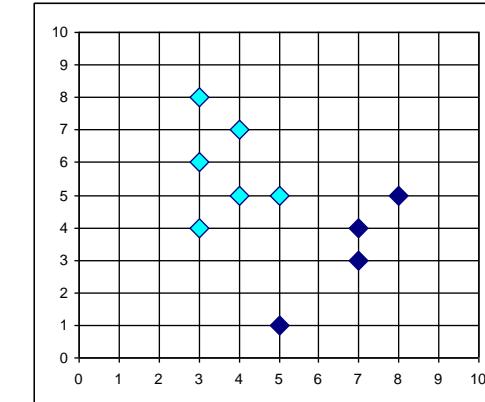


Update the cluster means

reassign



Update the cluster means



# K-means Clustering – Details

---

- Initial centroids are often chosen randomly
  - Example: Pick one point at random, then  $k-1$  other points, each as far away as possible from the previous points
- The centroid is (typically) the mean of the points in the cluster.
- ‘Nearest’ is measured by Euclidean distance, cosine similarity, correlation, etc.
- Most of the convergence happens in the first few iterations.
  - Often the stopping condition is changed to ‘Until relatively few points change clusters’
- Complexity is  $\mathcal{O}(t k n)$ 

$n$  is # objects,  $k$  is # clusters, and  $t$  is # iterations.

# Comments on the *K-Means* Method

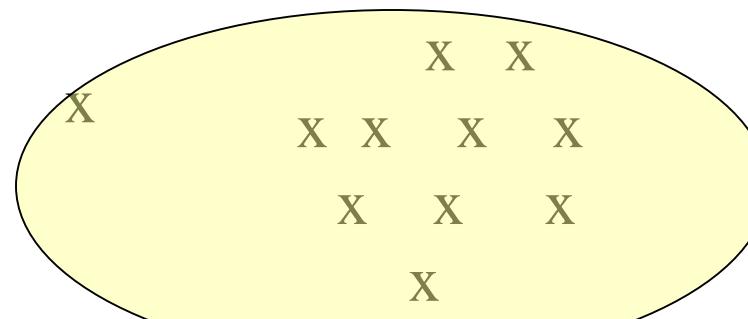
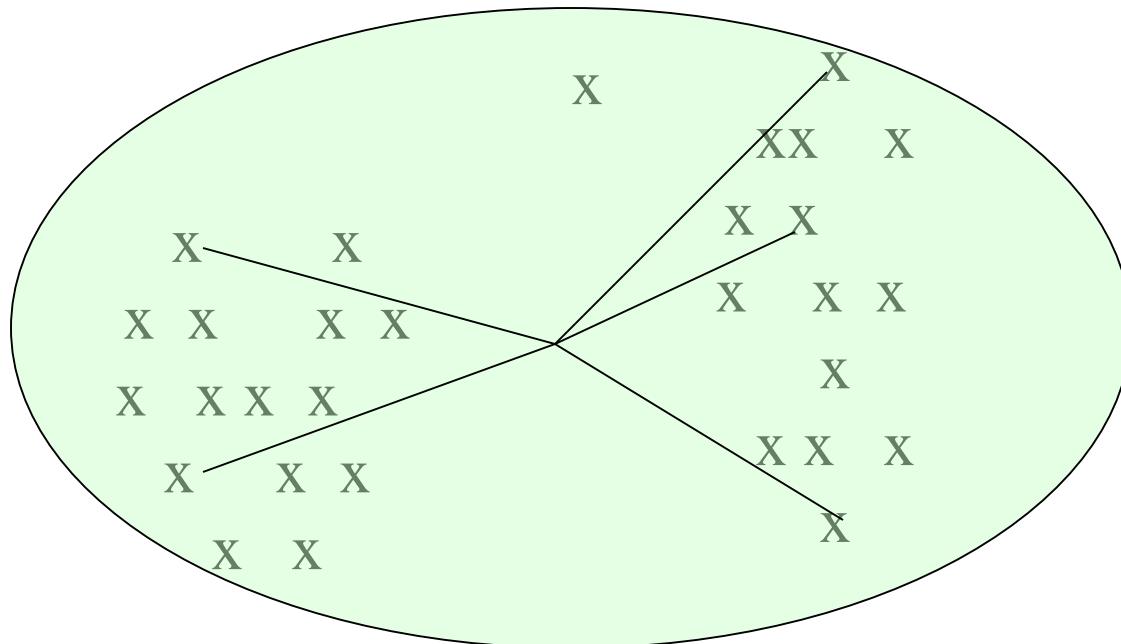
---

- Strength
  - Simple and works well for “regular” disjoint clusters
  - Relatively efficient and scalable (normally,  $k, t \ll n$ )
- Weakness
  - Need to specify  $k$ , the *number* of clusters, in advance
  - Depending on initial centroids, may terminate at a *local optimum*
  - Sensitive to noisy data and *outliers*
  - Not suitable for clusters of
    - Different sizes
    - Non-convex shapes

# Example: Picking $k$

---

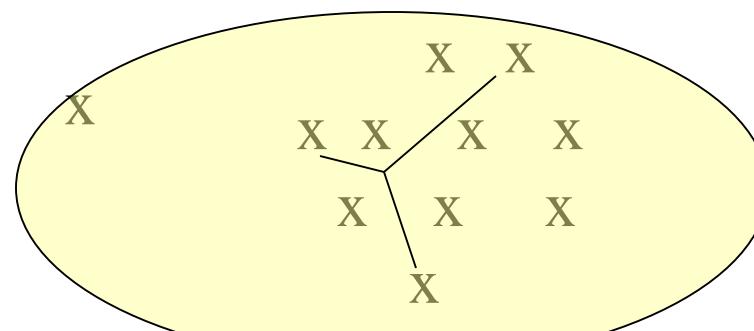
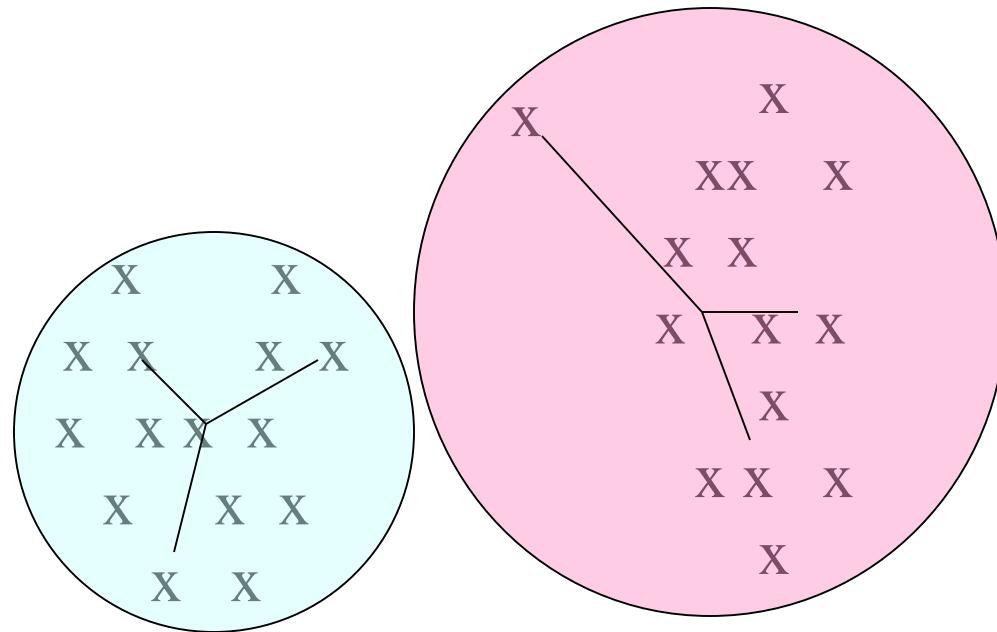
**Too few;**  
many long  
distances  
to centroid.



# Example: Picking $k$

---

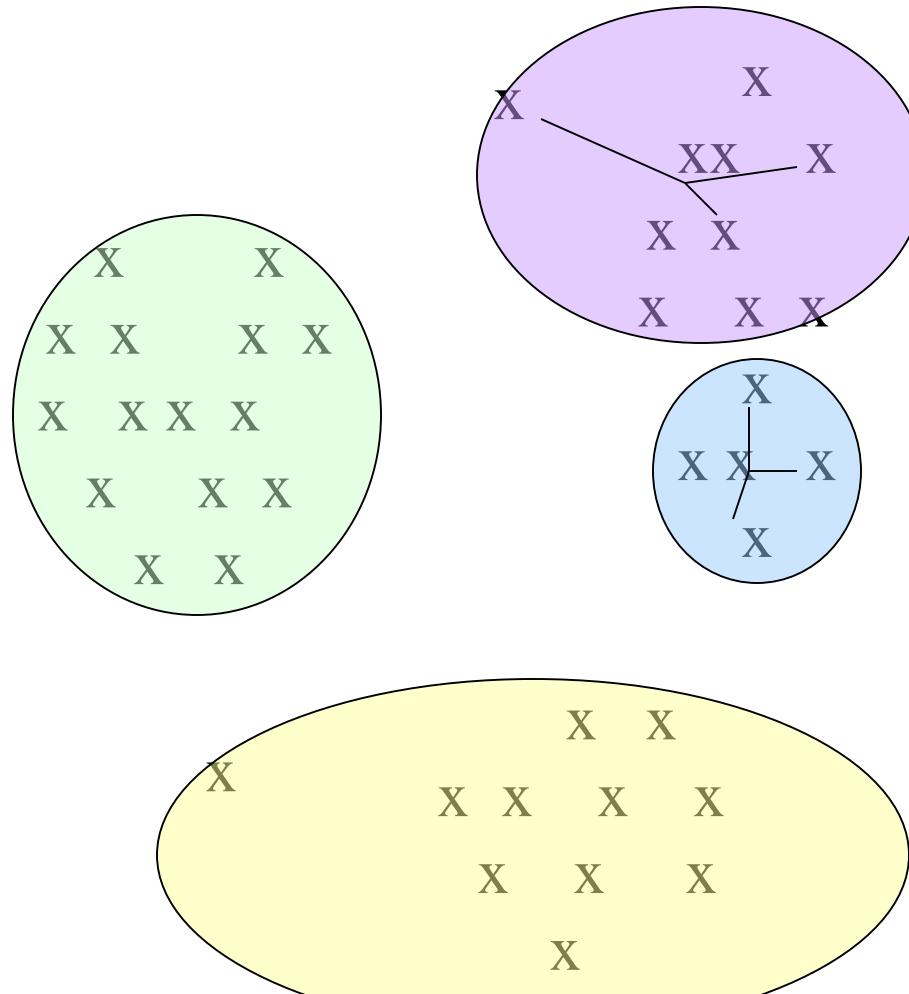
Just right;  
distances  
rather short.



# Example: Picking $k$

---

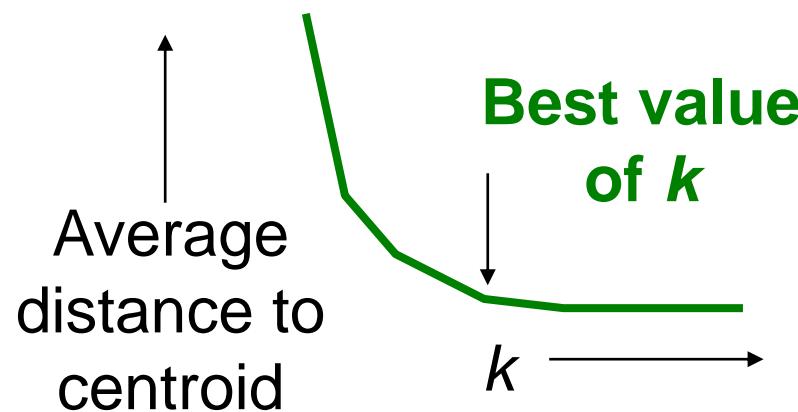
Too many;  
little improvement  
in average  
distance.



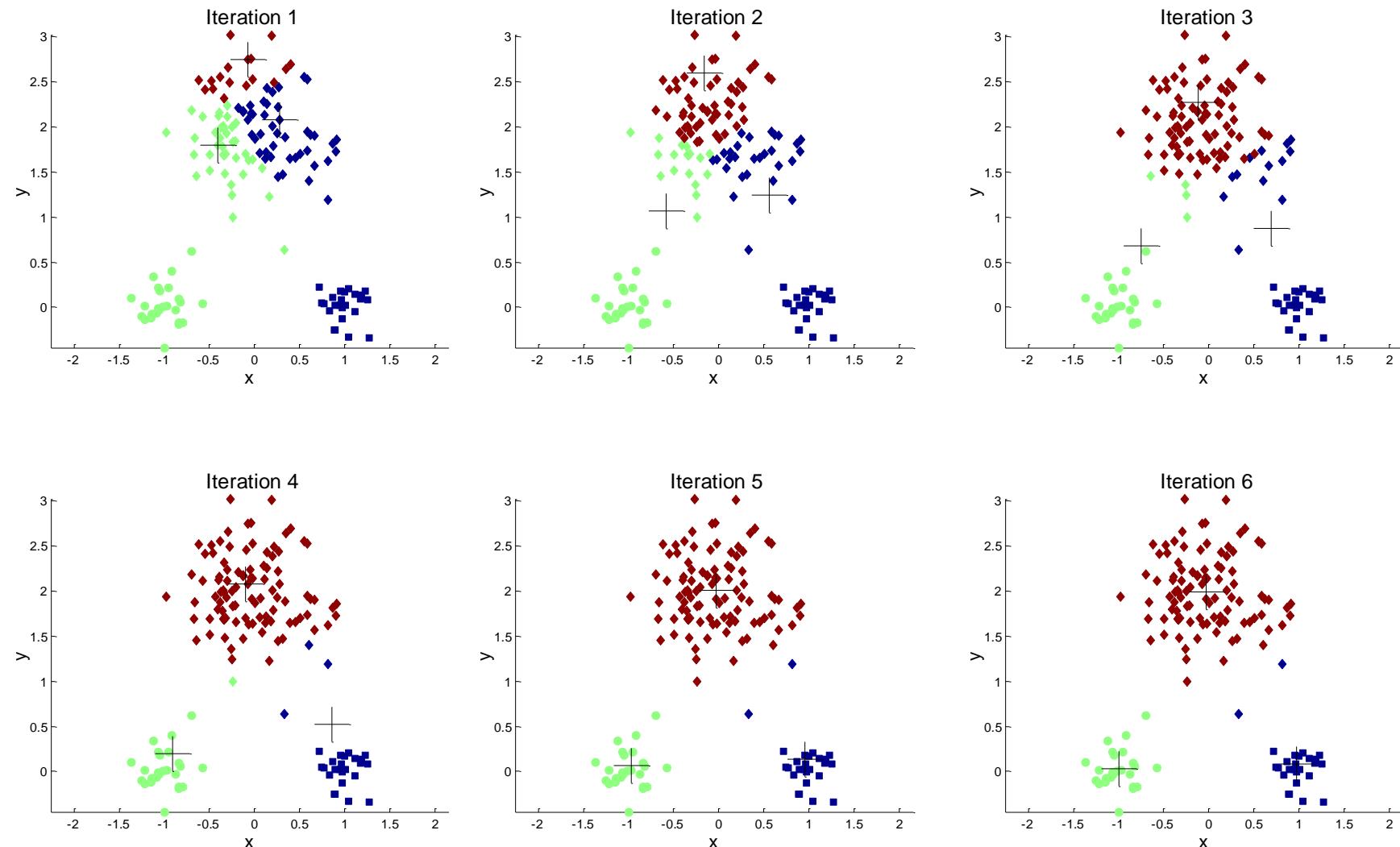
# Getting the $k$ right

## How to select $k$ ?

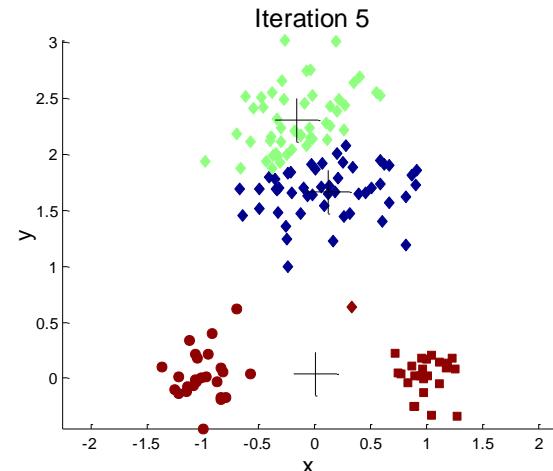
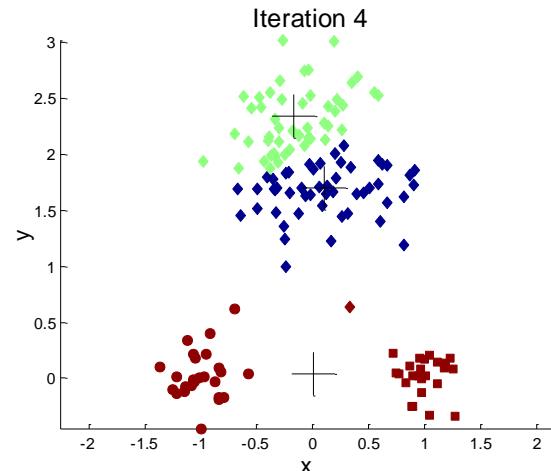
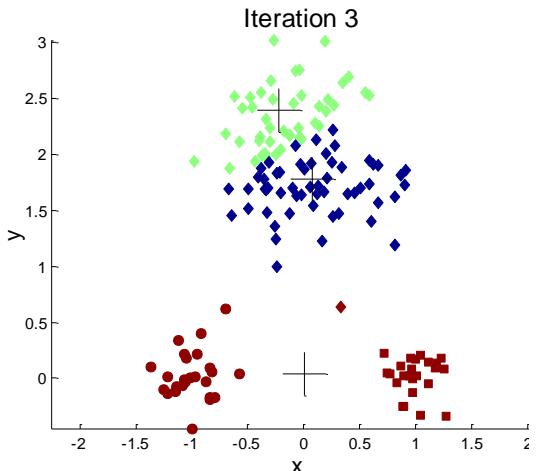
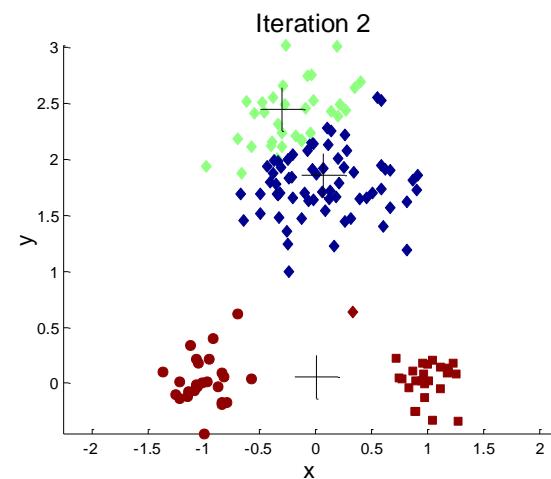
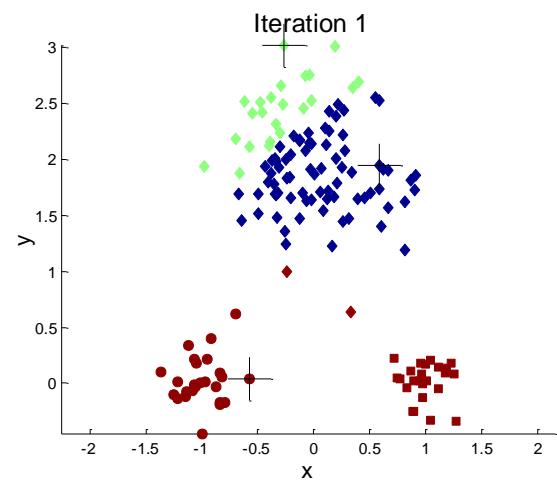
- Try different  $k$ , looking at the change in the average distance to centroid (or SSE) as  $k$  increases
- Average falls rapidly until right  $k$ , then changes little



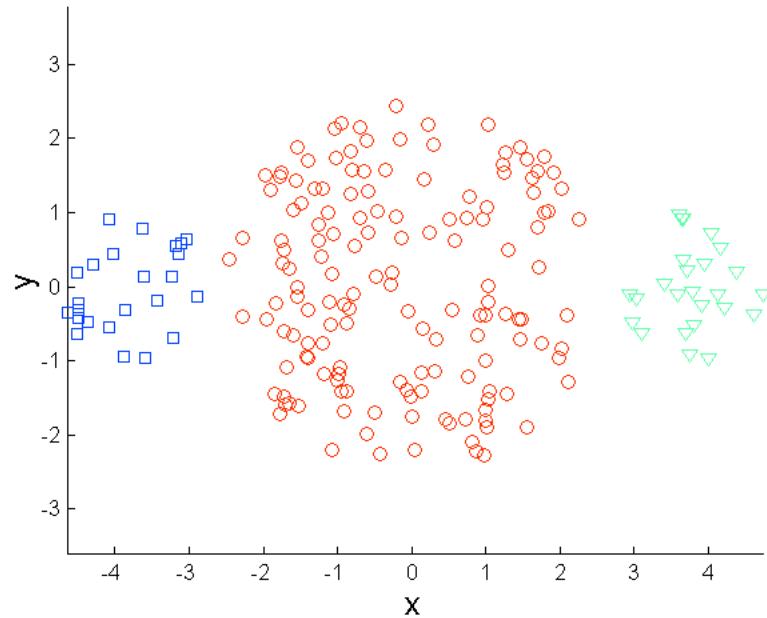
# Importance of Choosing Initial Centroids – Case 1



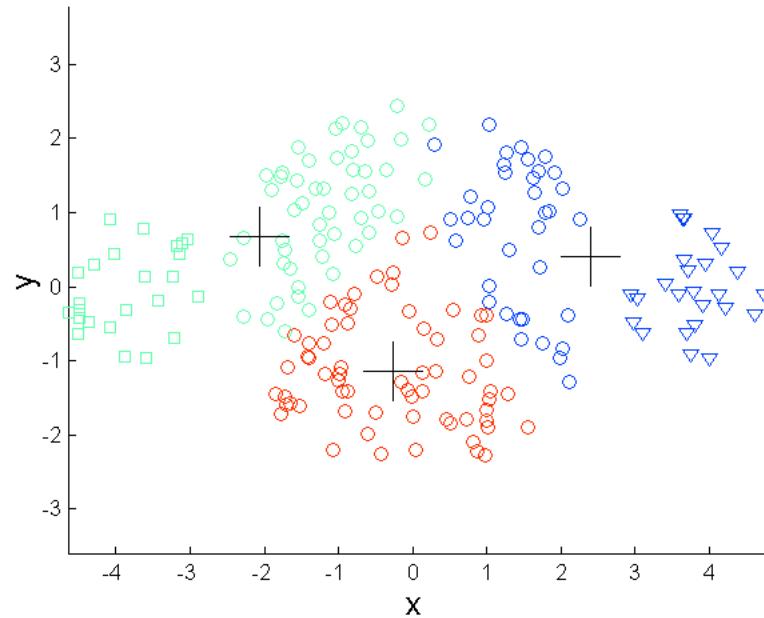
# Importance of Choosing Initial Centroids – Case 2



# Limitations of K-means: Differing Sizes

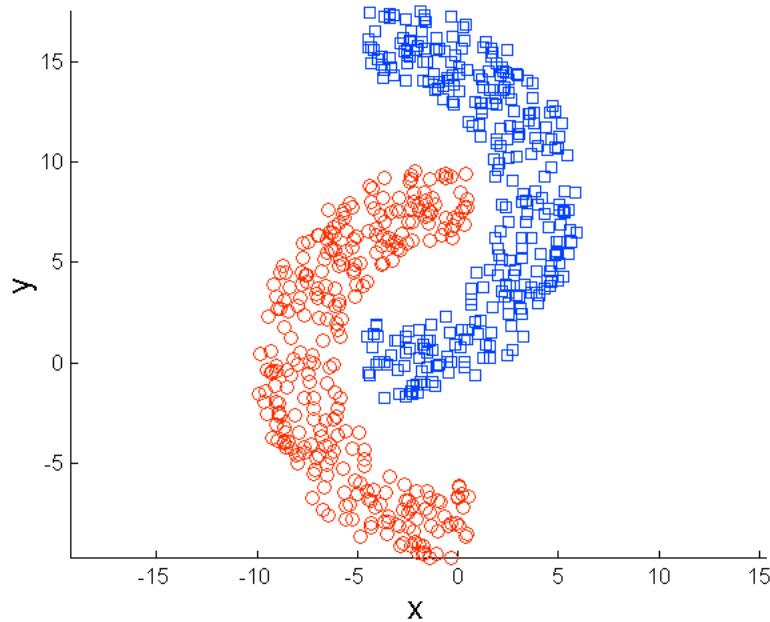


Original Points

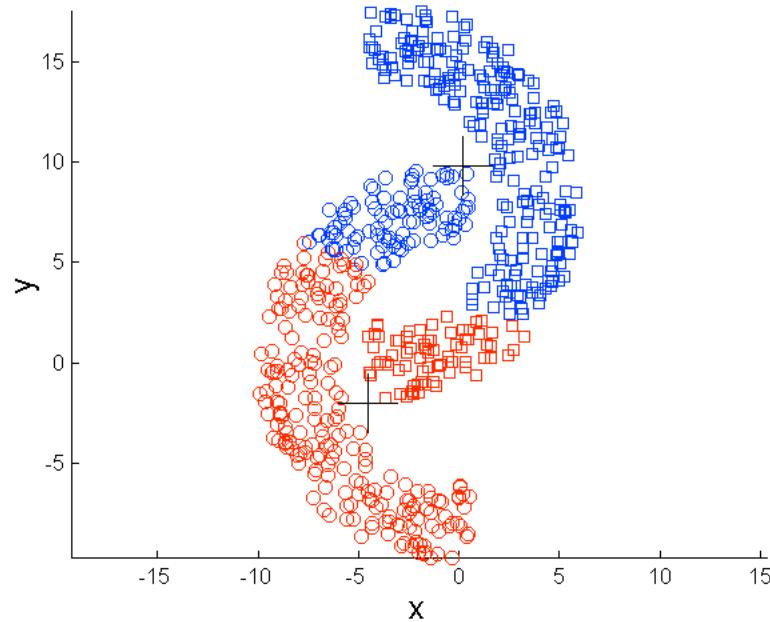


K-means (3 Clusters)

# Limitations of K-means: Non-convex Shapes

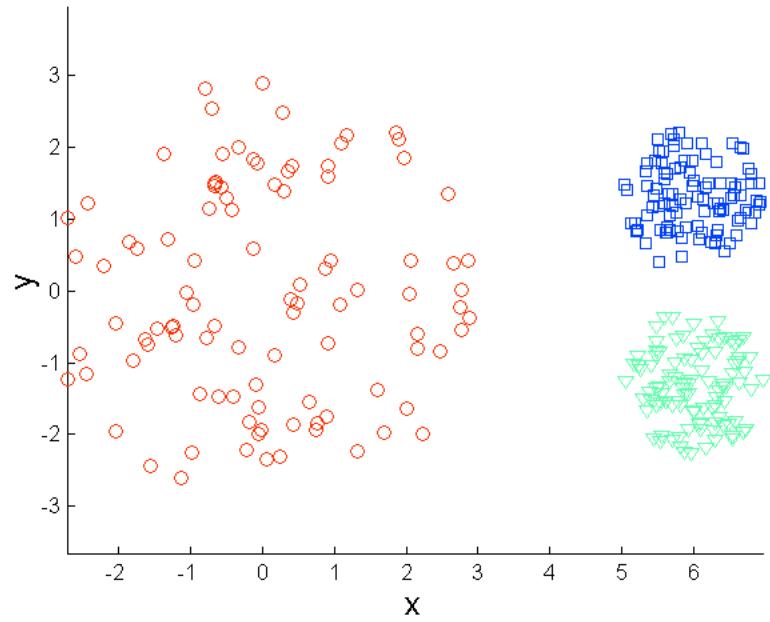


Original Points

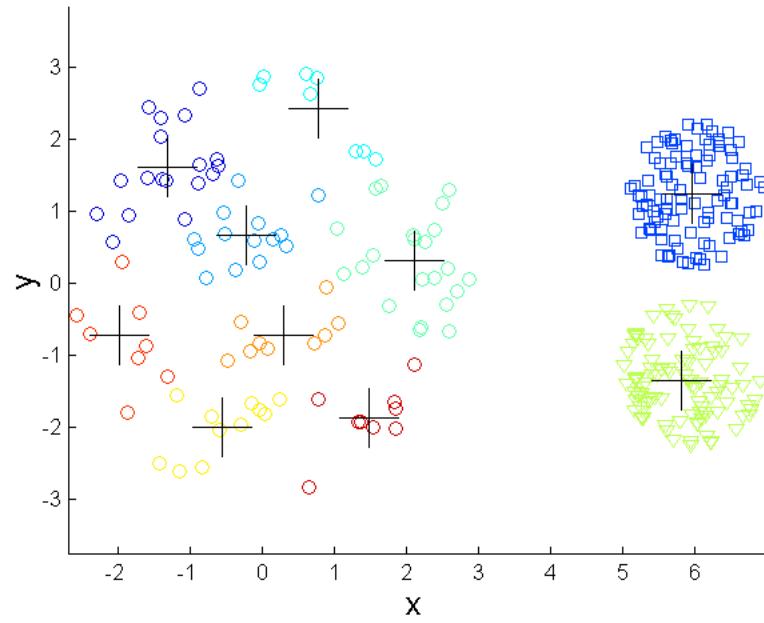


K-means (2 Clusters)

# Overcoming K-means Limitations

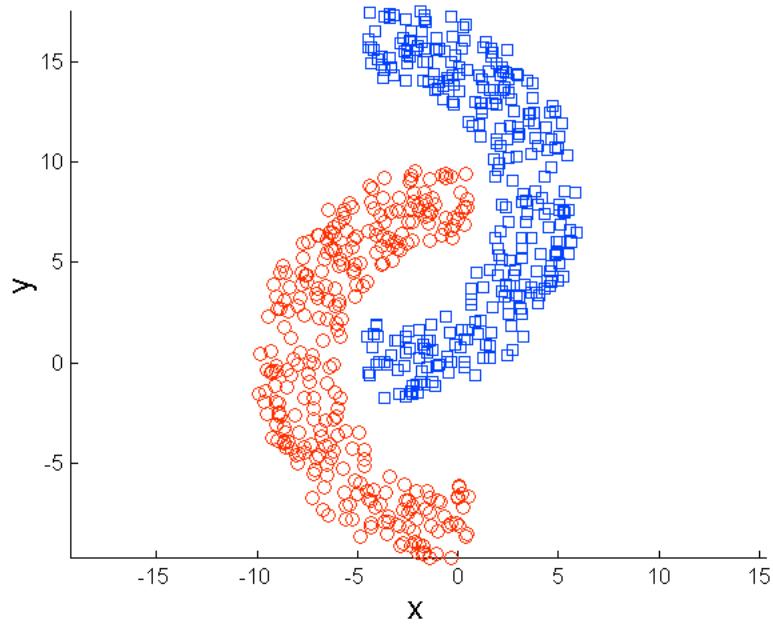


Original Points

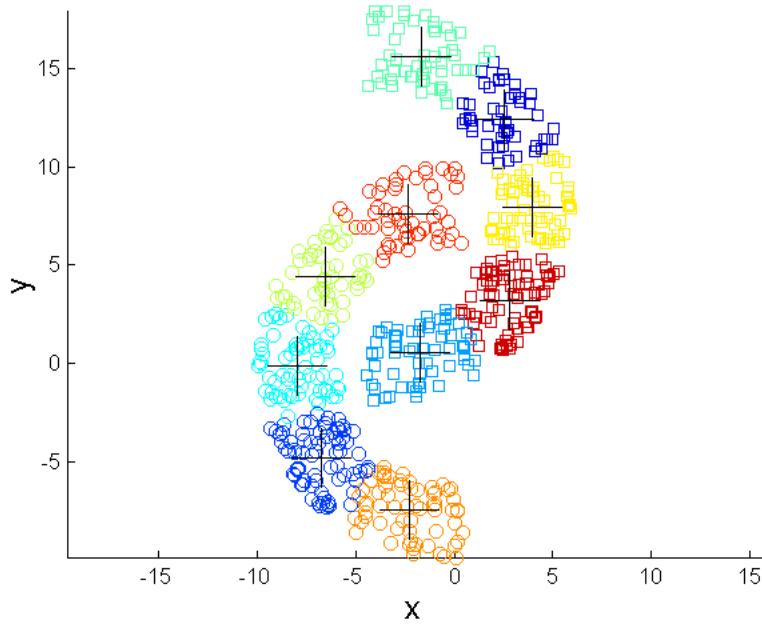


K-means Clusters

# Overcoming K-means Limitations



Original Points



K-means Clusters

# Assignment 3

---

- Implement k-means clustering
- Evaluate the results

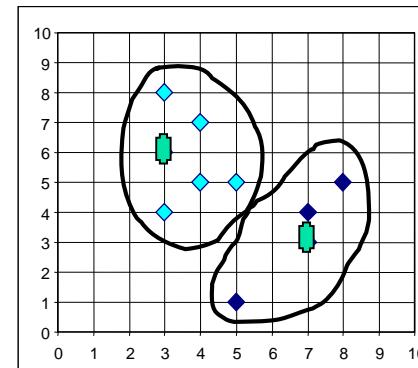
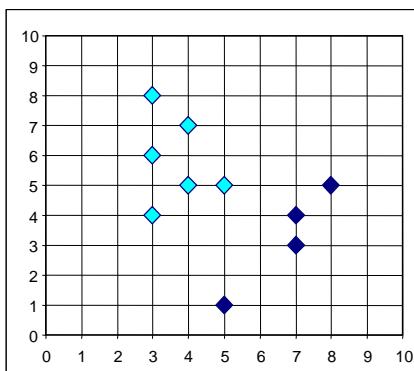
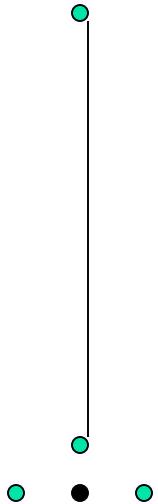
# Variations of the *K-Means* Method

---

- 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
  - A mixture of categorical and numerical data: *k-prototype* method

# K-Medoids Method

- The k-means algorithm is sensitive to outliers !
  - Since an object with an extremely large value may substantially distort the **mean** of the data.
- K-Medoids: Instead of using the **mean** as cluster representative, use **medoid**, the **most centrally located** object in a cluster.
- Possible number of solutions?



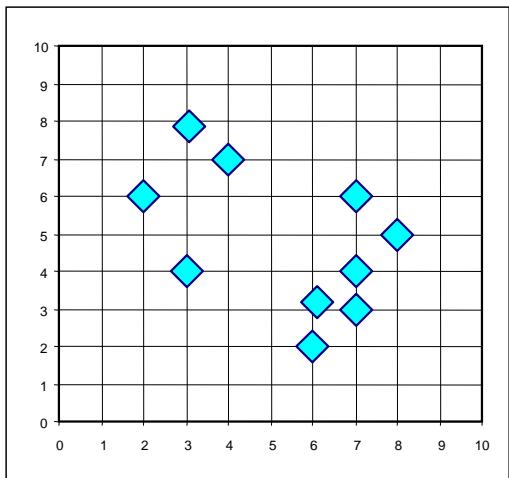
# The *K-Medoids* Clustering Method

---

PAM (Partitioning Around Medoids) (Kaufman and Rousseeuw, 1987)

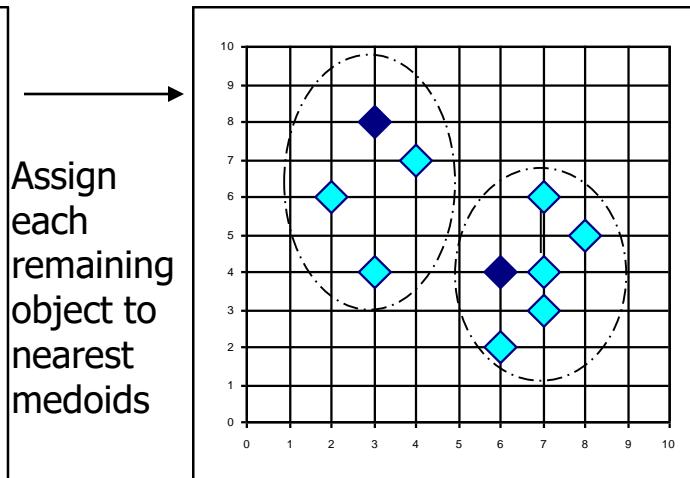
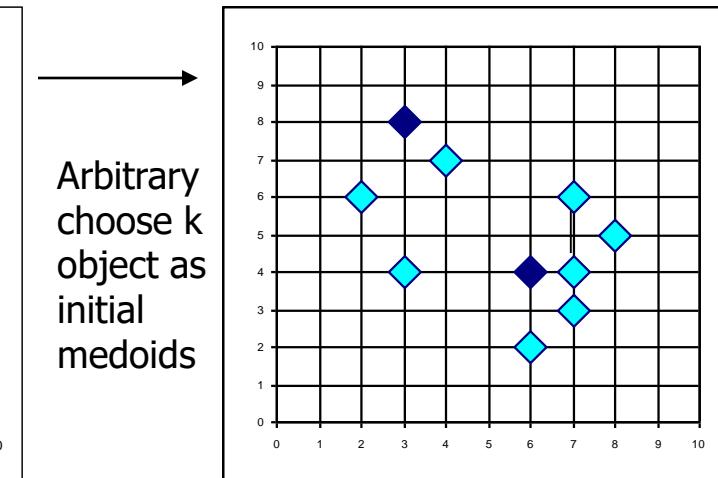
- Arbitrarily select  $k$  objects as medoid
- Assign each data object in the given data set to most similar medoid.
- For each nonmedoid object  $O'$  and medoid object  $O$ 
  - Compute total cost,  $S$ , of swapping the **medoid** object  $O$  to  $O'$  (cost as total sum of absolute error)
- If  $\min S < 0$ , then swap  $O$  with  $O'$
- Repeat until there is no change in the medoids.

# A Typical K-Medoids Algorithm (PAM)

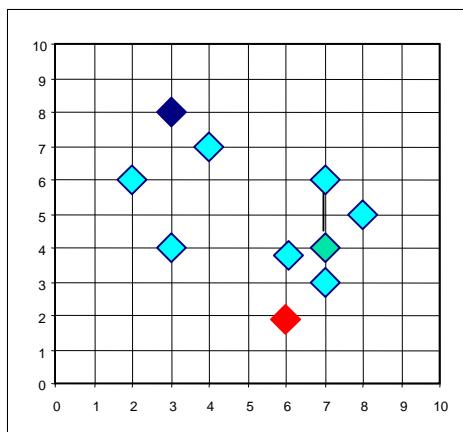


K=2

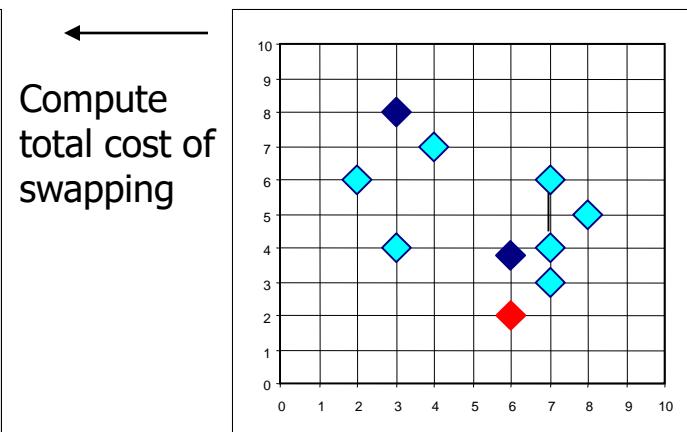
**Do loop**  
**Until no change**



Select a nonmedoid object,  $O_{random}$



Swapping  $O$  and  $O_{random}$   
If quality is improved.



# What Is the Problem with PAM?

---

- Pam is more robust than k-means in the presence of noise and outliers
- Pam works efficiently for small data sets but does not **scale well** for large data sets.
  - Complexity?  
 $n$  is # of data,  $k$  is # of clusters

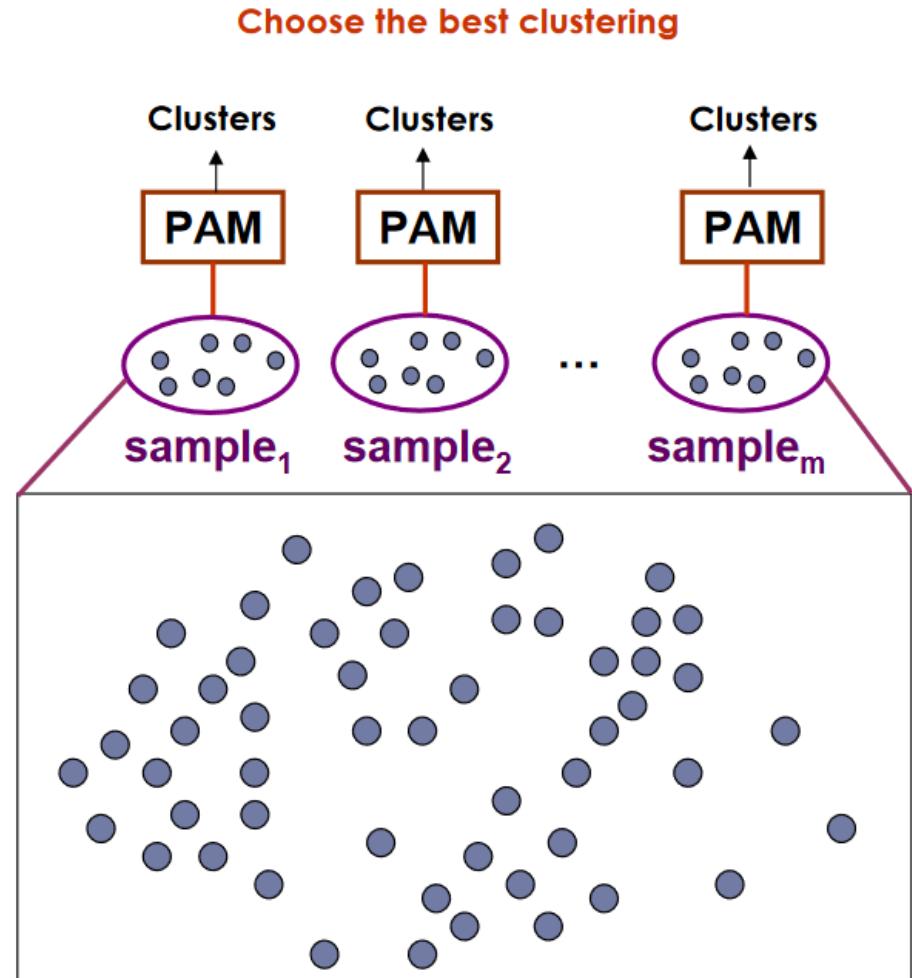
# What Is the Problem with PAM?

---

- Pam is more robust than k-means in the presence of noise and outliers
- Pam works efficiently for small data sets but does not **scale well** for large data sets.
  - Complexity?  $O(k(n-k)^2)$   
 $n$  is # of data,  $k$  is # of clusters

# CLARA (Clustering Large Applications) (1990)

- CLARA (Kaufmann and Rousseeuw in 1990)
- Draws *multiple samples* of the data set, applies PAM on each sample, and gives the best clustering as the output



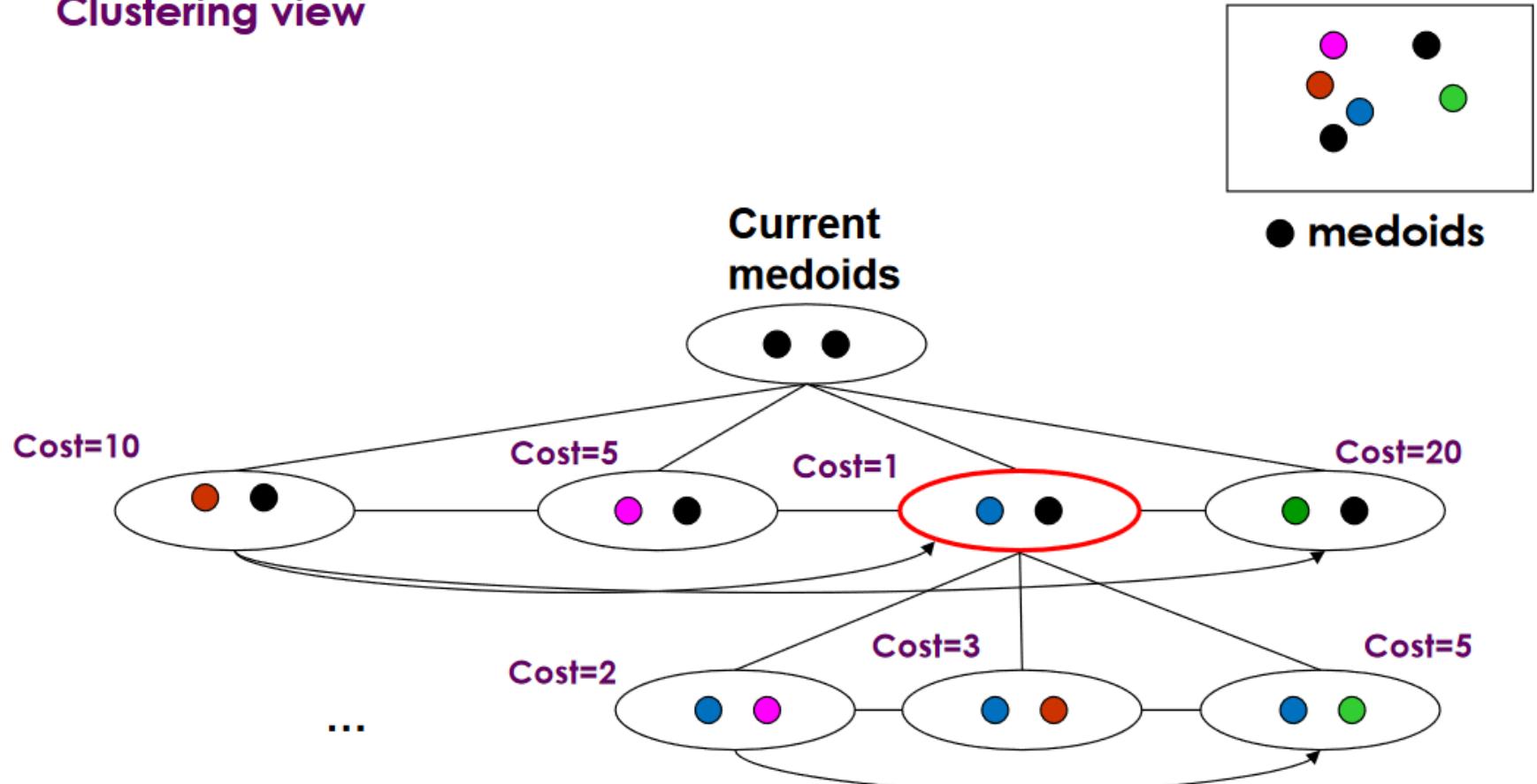
# *CLARANS ("Randomized" CLARA) (1994)*

---

- *CLARANS* (A Clustering Algorithm based on Randomized Search) (Ng and Han'94)
- The clustering process can be represented as searching a graph where every node is a potential solution, that is, a set of  $k$  medoids

# Search graph

## Clustering view



# *CLARANS* ("Randomized" CLARA) (1994)

---

- *CLARANS* (A Clustering Algorithm based on Randomized Search) (Ng and Han'94)
- The clustering process can be represented as searching a graph where every node is a potential solution, that is, a set of  $k$  medoids
  - PAM examines all neighbors for local minimum
  - CLARA works on subgraphs of samples
  - CLARANS examines neighbors dynamically
    - Limit the neighbors to explore (*maxneighbor*)
    - If local optimum is found, start with new randomly selected node in search for a new local optimum (*numlocal*)

# Cluster Analysis: Basic Concepts and Methods

---

- Cluster Analysis: Basic Concepts
- Similarity and distances
- Partitioning Methods
- Hierarchical Methods
- Density-Based Methods
- Probabilistic Methods
- Evaluation of Clustering

# Cluster Evaluation

---

- Determine clustering tendency of data, i.e. distinguish whether non-random structure exists
- Determine correct number of clusters
- Evaluate the cohesion and separation of the clustering without external information
- Evaluate how well the cluster results are compared to externally known results
- Compare different clustering algorithms/results

# Measures

---

- **Unsupervised (internal):** Used to measure the goodness of a clustering structure *without* respect to external information.
  - Sum of Squared Error (SSE)
- **Supervised (external):** Used to measure the extent to which cluster labels match externally supplied class labels.
  - Entropy
- **Relative:** Used to compare two different clustering results
  - Often an external or internal index is used for this function, e.g., SSE or entropy

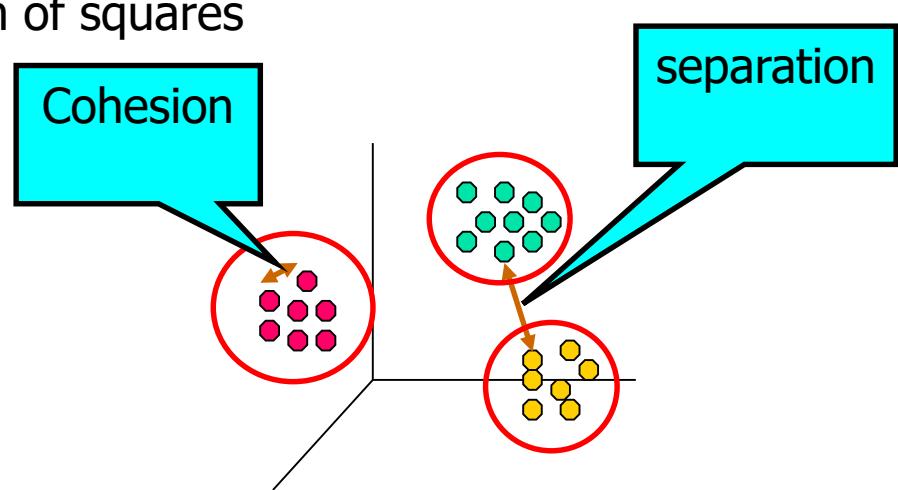
# Internal Measures: Cohesion and Separation

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

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

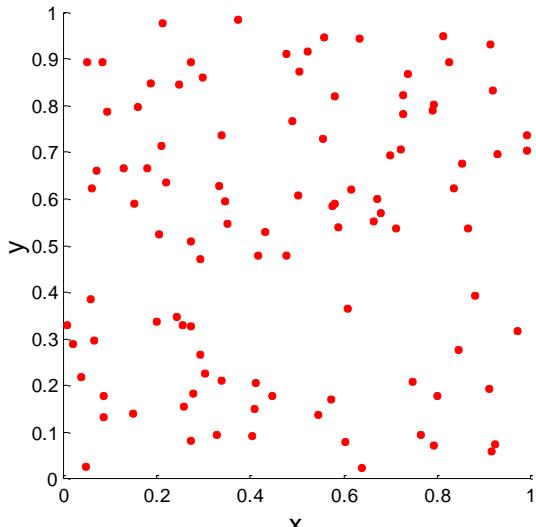
- Separation: between cluster sum of squares

$$BSS = \sum_i \sum_j (m_i - m_j)^2$$

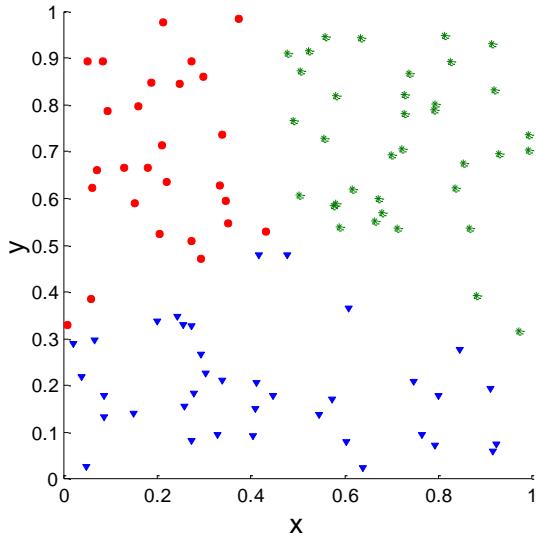


# Cluster Validity: Clusters found in Random Data

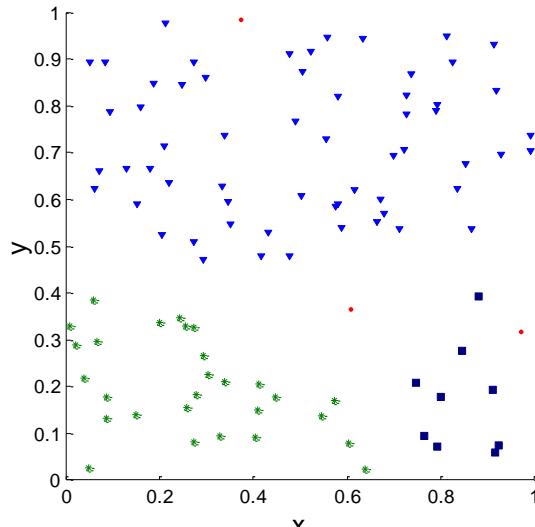
Random Points



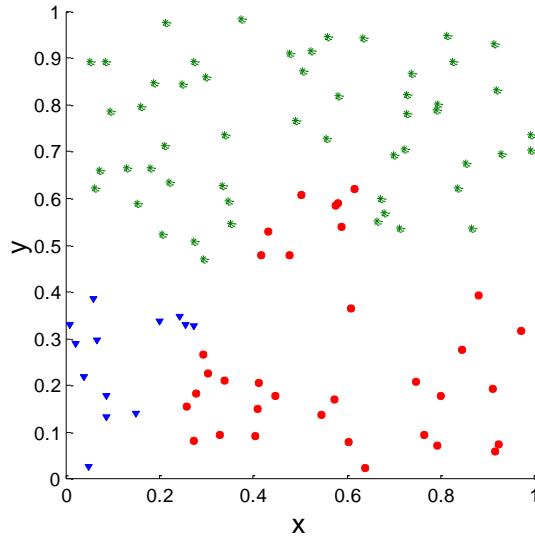
K-means



DBSCAN

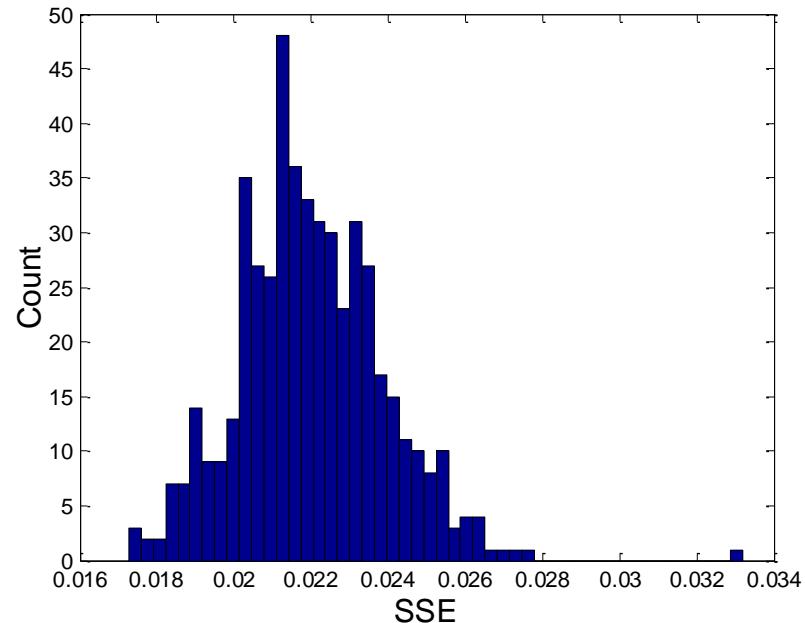
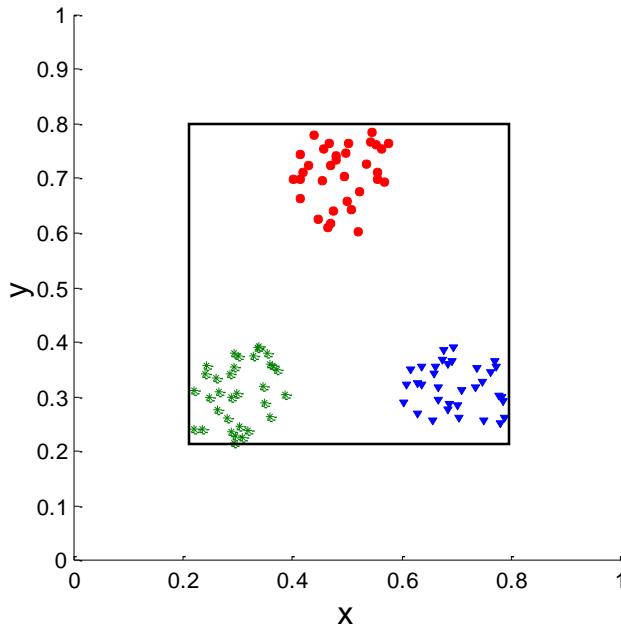


Complete Link



# Internal Measures: Cluster Validity

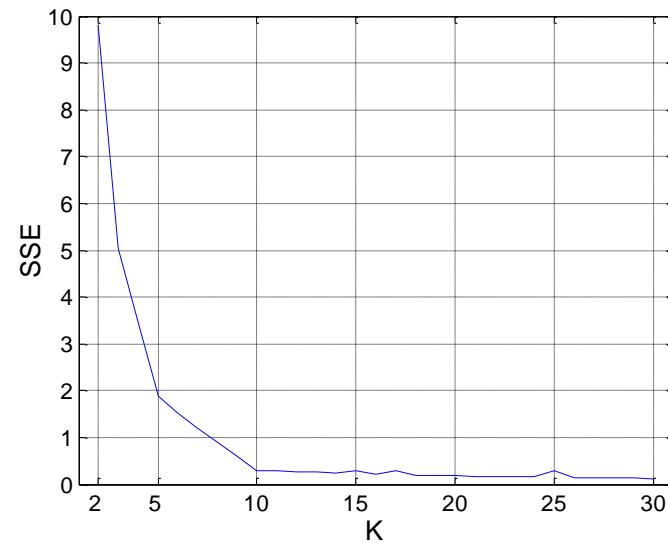
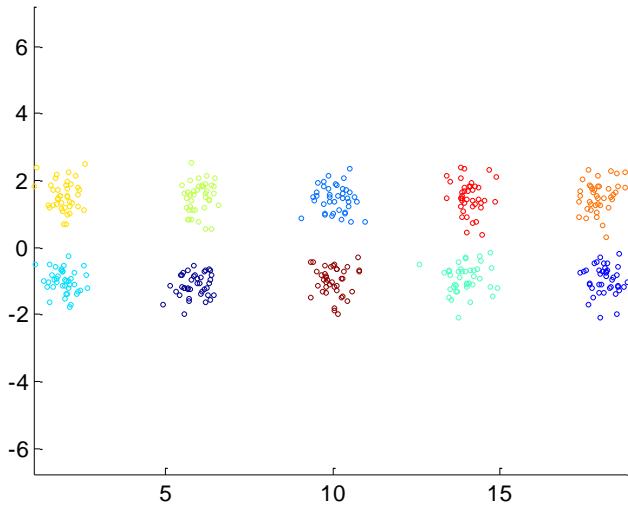
- Statistics framework for cluster validity
  - More “atypical” -> likely valid structure in the data
  - Use values resulting from random data as baseline
- Example
  - Clustering: SSE = 0.005
  - SSE of three clusters in 500 sets of random data points



# Internal Measures: number of clusters

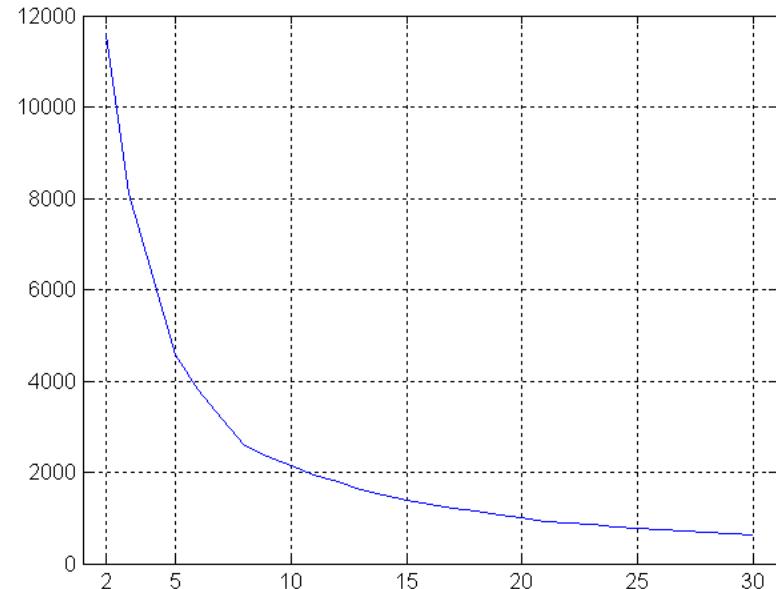
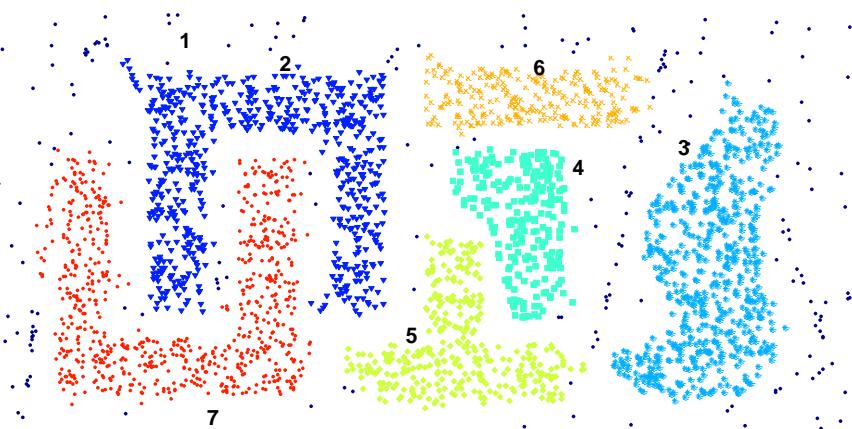
---

- Good for comparing two clusterings
- Can also be used to estimate the number of clusters
  - Elbow method: use turning point in the curve of SSE wrt # of clusters



# Internal Measures: Number of clusters

- Another example of a more complicated data set with varying number of clusters



SSE of clusters found using K-means

# External Measures

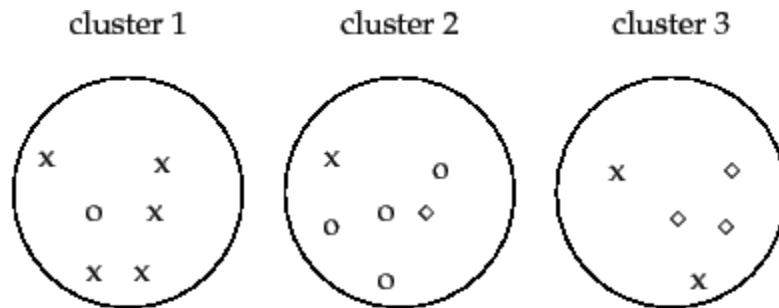
---

- Compare cluster results with “ground truth” or manually clustering
- Still different from classification measures
- Classification-oriented measures: entropy/purity based, precision and recall based
- Similarity-oriented measures: Jaccard scores

# External Measures: Classification-Oriented Measures

---

- Entropy based measures: the degree to which each *cluster* consists of objects of a single class
- Purity: based on majority class in each cluster



► **Figure 16.1** Purity as an external evaluation criterion for cluster quality. Majority class and number of members of the majority class for the three clusters are: x, 5 (cluster 1); o, 4 (cluster 2); and  $\diamond$ , 3 (cluster 3). Purity is  $(1/17) \times (5 + 4 + 3) \approx 0.71$ .

## External Measures: Classification-Oriented Measures

---

- BCubed Precision and recall: measures precision and recall associated with each *object*
  - Precision of an object: proportion of objects in the same cluster belong to the same category
  - Recall of an object: proportion of objects of the same category are assigned to the same cluster
  - Bcubed precision and recall are the average precision and recall of all objects

# BCubed precision and recall

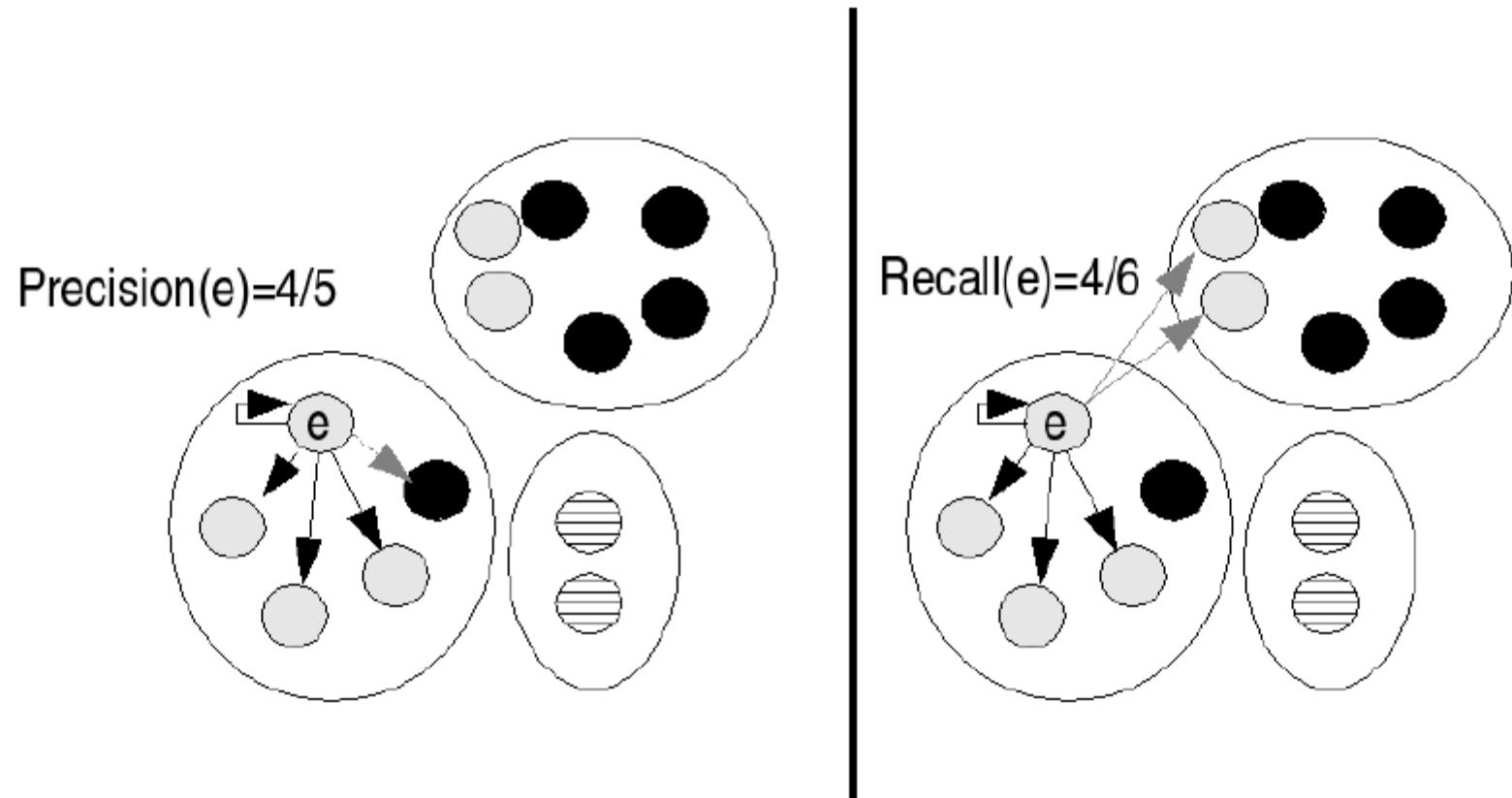


Figure 10: Example of computing the BCubed precision and recall for one item

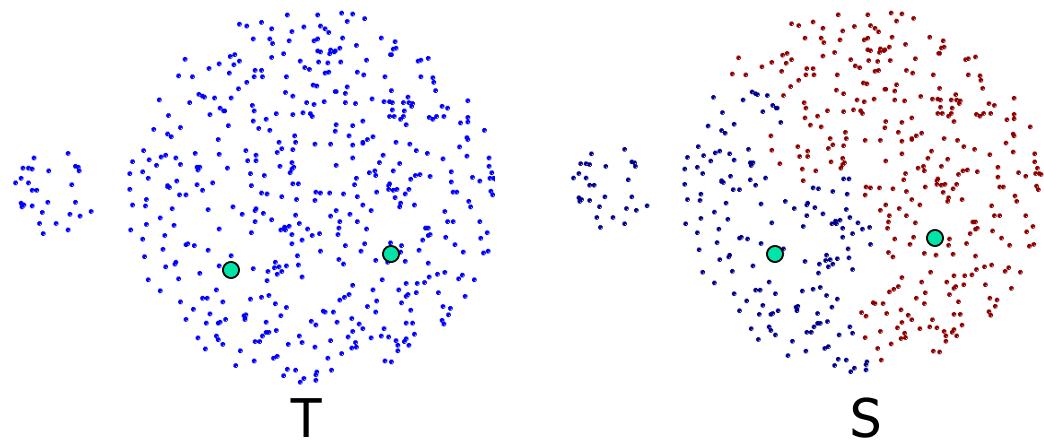
# External Measure: Similarity-Oriented Measures

Given a reference clustering T and clustering S

- $f_{00}$ : number of pair of points belonging to different clusters in both T and S
- $f_{01}$ : number of pair of points belonging to different cluster in T but same cluster in S
- $f_{10}$ : number of pair of points belonging to same cluster in T but different cluster in S
- $f_{11}$ : number of pair of points belonging to same clusters in both T and S

$$Rand = \frac{f_{00} + f_{11}}{f_{00} + f_{01} + f_{10} + f_{11}}$$

$$Jaccard = \frac{f_{11}}{f_{01} + f_{10} + f_{11}}$$

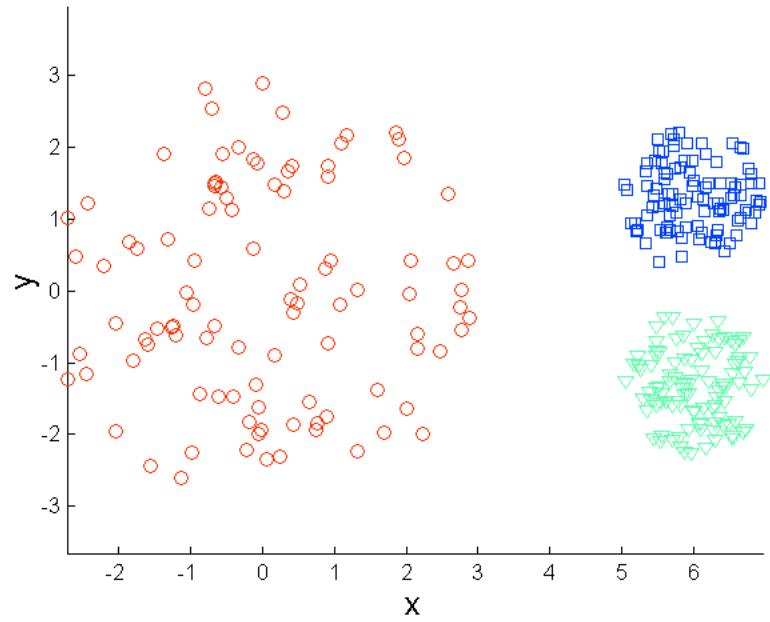


# Cluster Analysis: Basic Concepts and Methods

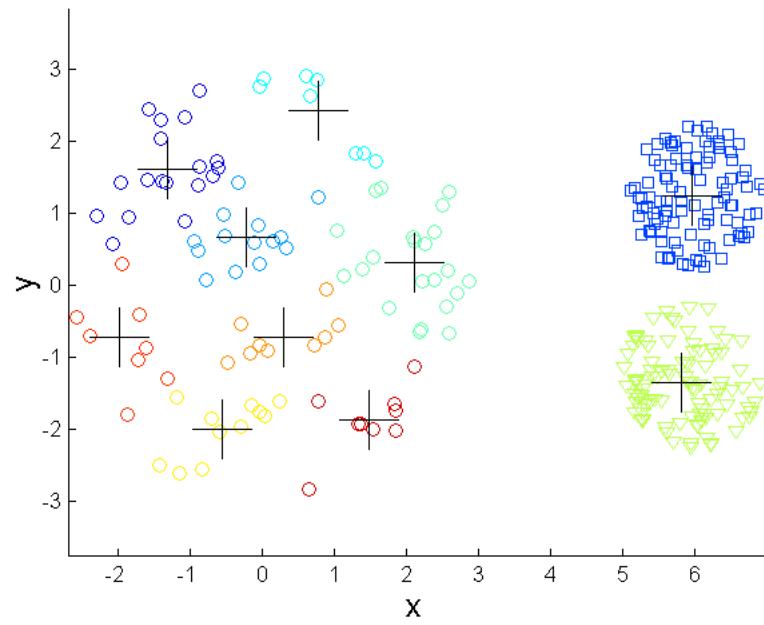
---

- Cluster Analysis: Basic Concepts
- Similarity and distances
- Partitioning Methods
- Hierarchical Methods
- Density-Based Methods
- Probabilistic Methods
- Evaluation of Clustering

# Overcoming K-means Limitations

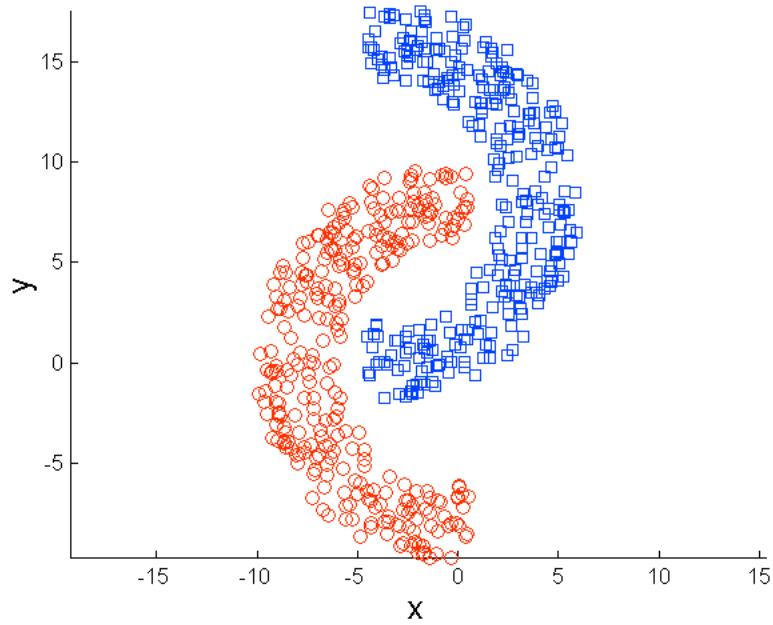


Original Points

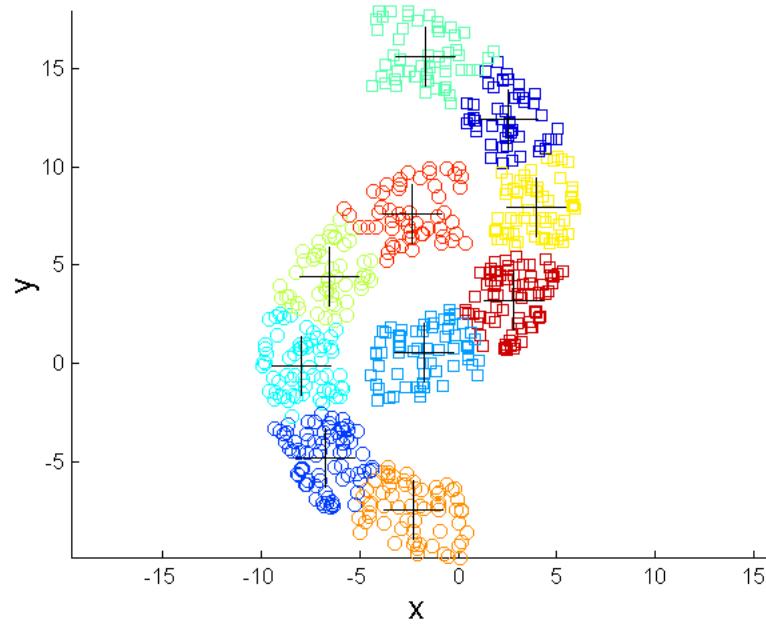


K-means Clusters

# Overcoming K-means Limitations



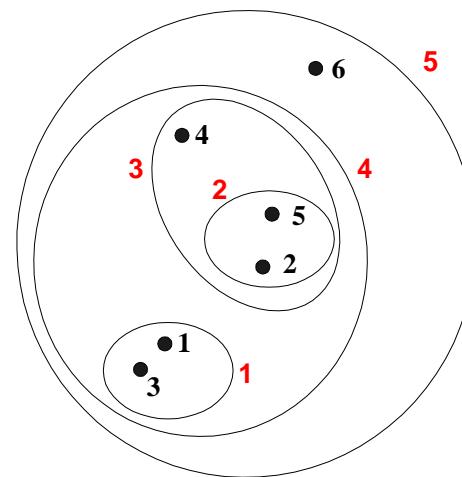
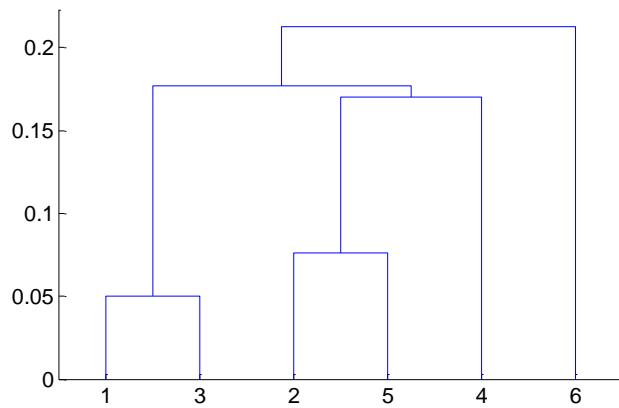
Original Points

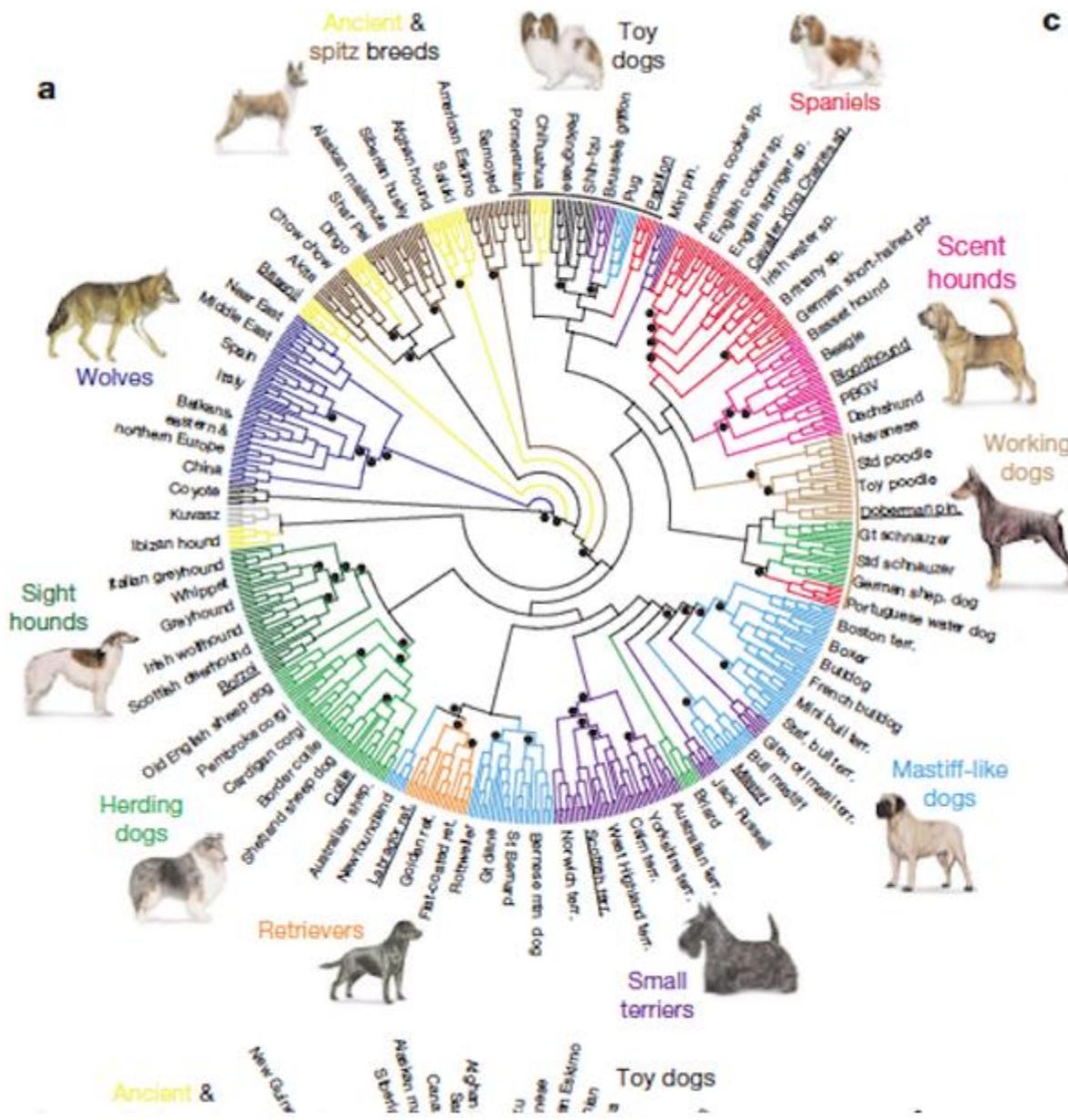


K-means Clusters

# Hierarchical Clustering

- Produces a set of nested clusters
- Can be visualized as a dendrogram, a tree like diagram
  - Y-axis measures closeness
  - Clustering obtained by cutting at desired level
- Do not have to assume any particular number of clusters
- May correspond to meaningful taxonomies





# Hierarchical Clustering

---

- Two main types of hierarchical clustering
  - Agglomerative (AGNES)
    - 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 (DIANA)
    - Start with one, all-inclusive cluster
    - At each step, split a cluster until each cluster contains a point (or there are  $k$  clusters)

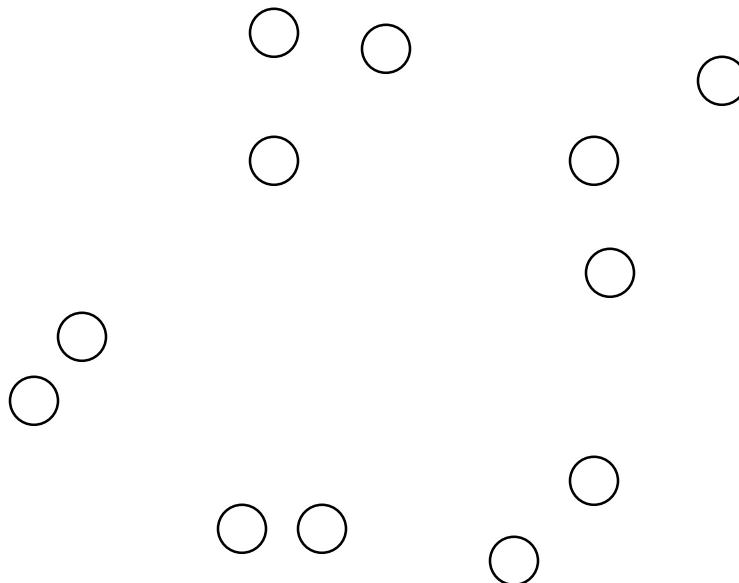
# Agglomerative Clustering Algorithm

---

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

# Starting Situation

- Start with clusters of individual points and a proximity matrix



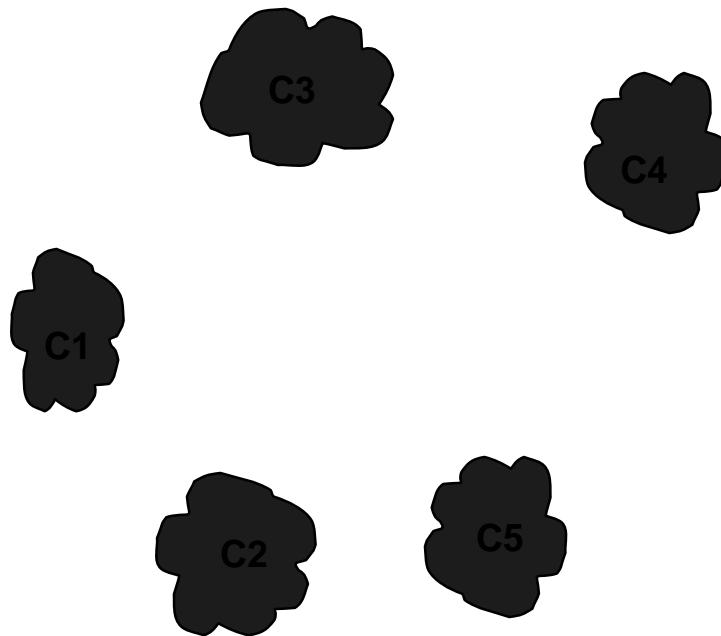
	p1	p2	p3	p4	p5	...
p1						
p2						
p3						
p4						
p5						
.						
.						

**Proximity Matrix**

p1   p2   p3   p4   ...   p9   p10   p11   p12

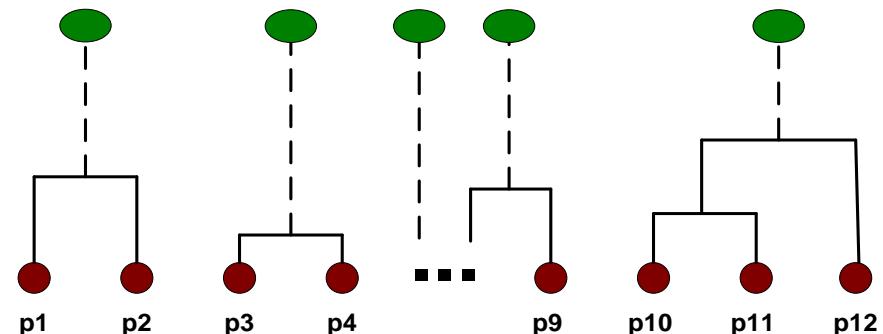
# Intermediate Situation

---



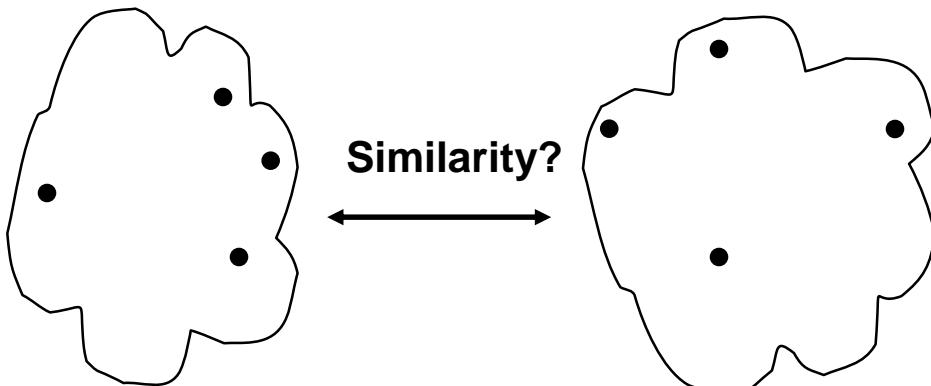
	C1	C2	C3	C4	C5
C1					
C2					
C3					
C4					
C5					

Proximity Matrix



# How to Define Inter-Cluster Similarity

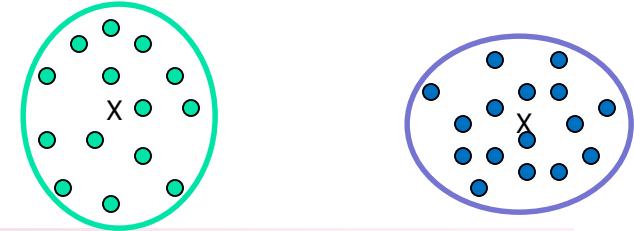
---



	p1	p2	p3	p4	p5	...
p1						
.						

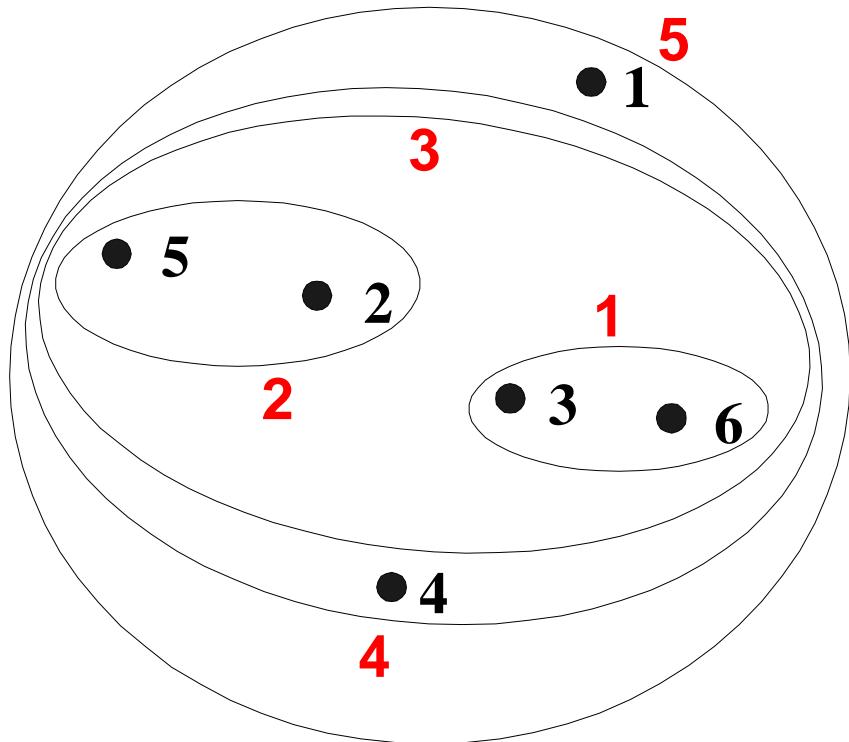
- **Proximity Matrix**

# Distance between Clusters

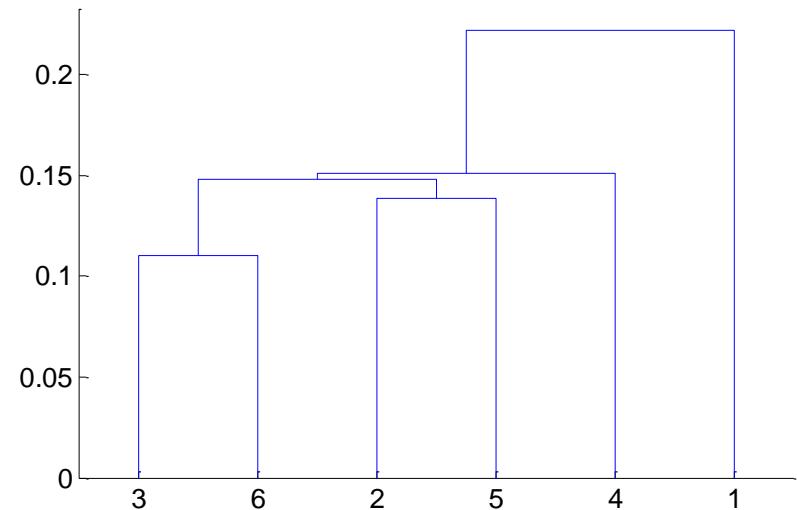


- Single link: smallest distance between an element in one cluster and an element in the other, i.e.,  $\text{dist}(K_i, K_j) = \min(t_{ip}, t_{jq})$
- Complete link: largest distance between an element in one cluster and an element in the other, i.e.,  $\text{dist}(K_i, K_j) = \max(t_{ip}, t_{jq})$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

# Hierarchical Clustering: MIN



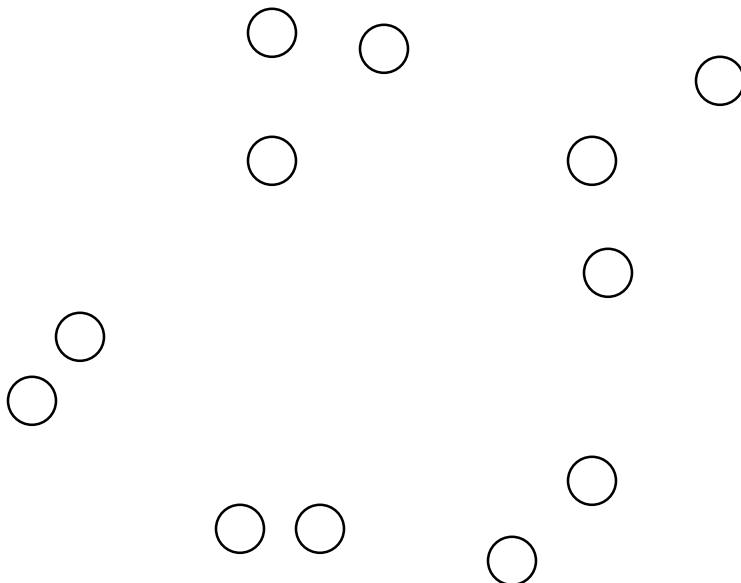
Nested Clusters



Dendrogram

# View points/similarities as a graph

- Start with clusters of individual points and a proximity matrix



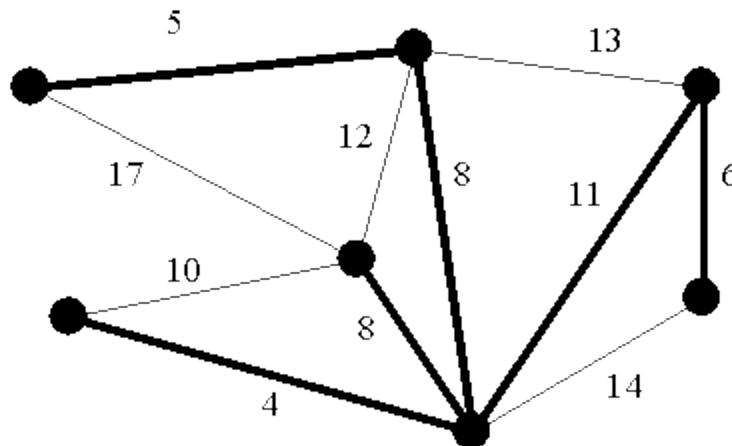
	p1	p2	p3	p4	p5	...
p1						
p2						
p3						
p4						
p5						
.						
.						

**Proximity Matrix**

p1   p2   p3   p4   ...   p9   p10   p11   p12

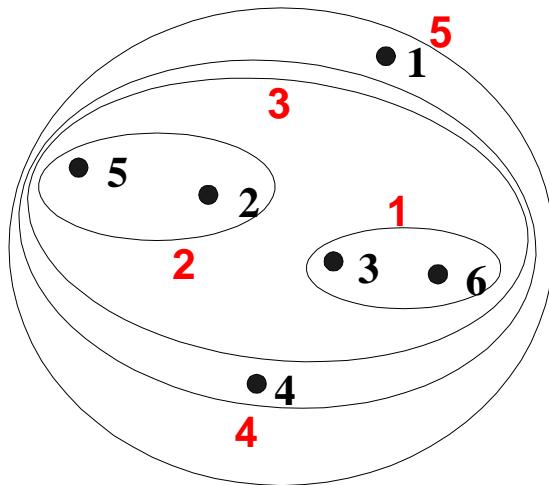
# Single link clustering and MST (Minimum Spanning Tree)

- An agglomerative algorithm using minimum distance (single-link clustering) essentially the same as Kruskal's algorithm for minimal spanning tree (MST)
- MST: a subgraph which is a tree and connects all vertices together that has the minimum weight
- Kruskal's algorithm: Add edges in increasing weight, skipping those whose addition would create a cycle
- Prim's algorithm: Grow a tree with any root node, adding the frontier edge with smallest weight

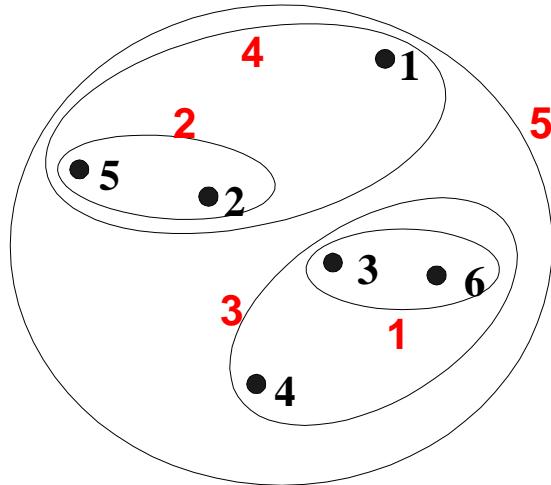


# Min vs. Max vs. Group Average

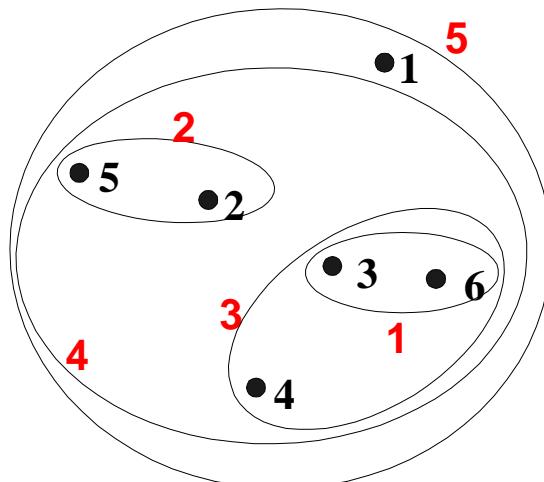
---



MIN



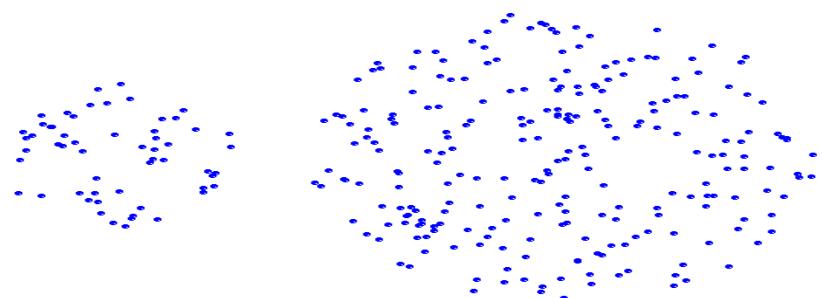
MAX



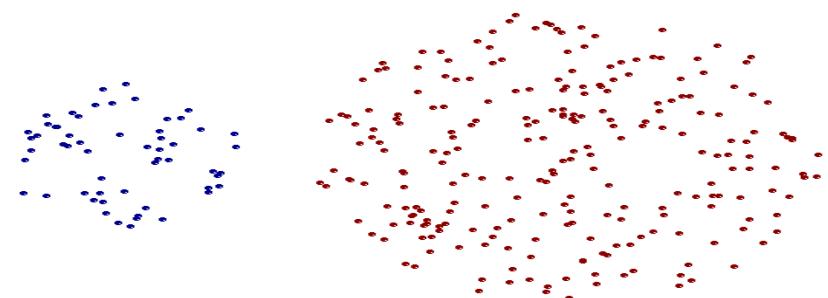
Group Average

# Strength of MIN

---



Original Points

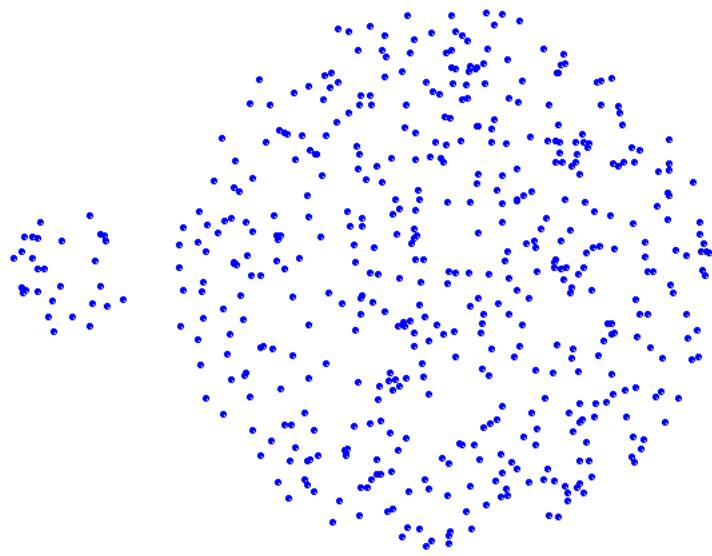


Two Clusters

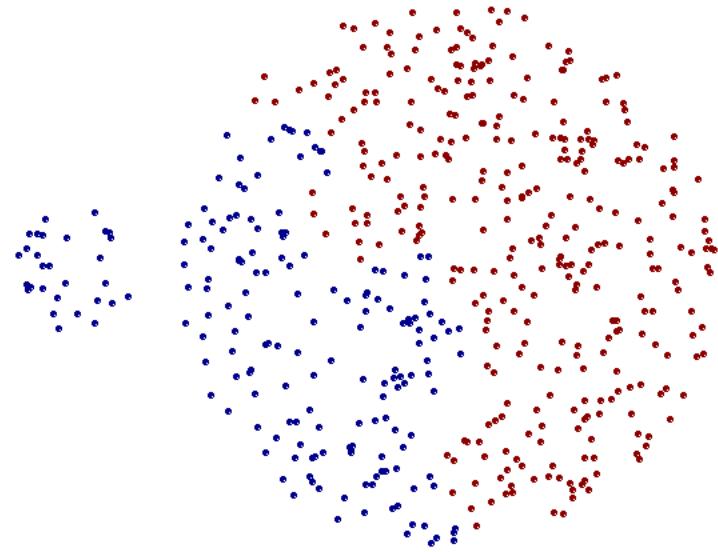
- Can handle clusters with varying sizes
- Can also handle non-elliptical shapes

# Limitations of MAX

---



**Original Points**

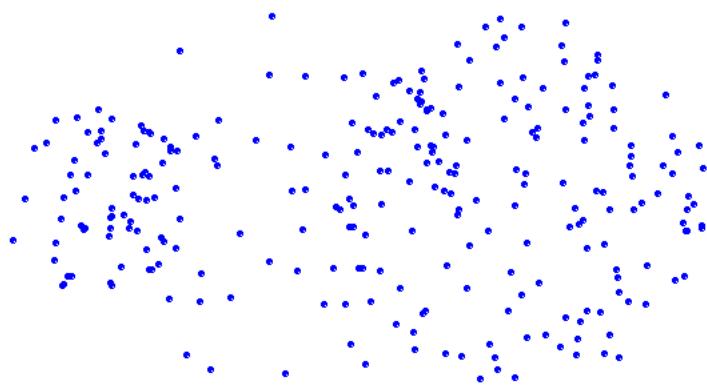


**Two Clusters**

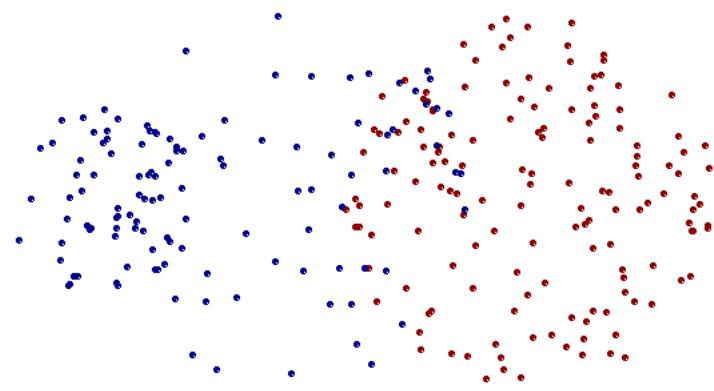
- Tends to break large clusters
- Biased towards globular clusters

# Limitations of MIN

---



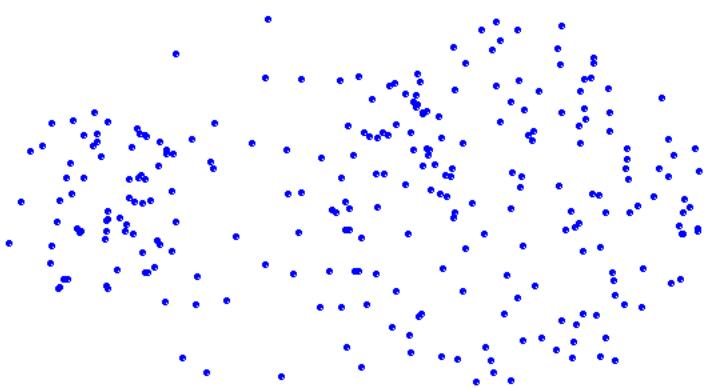
Original Points



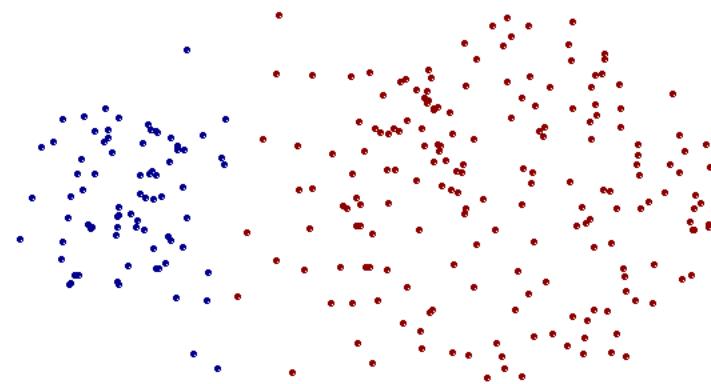
Two Clusters

- Chaining phenomenon
- Sensitive to noise and outliers

# Strength of MAX



Original Points



Two Clusters

- Less susceptible to noise and outliers

# Hierarchical Clustering: Group Average

---

- Compromise between Single and Complete Link
- Strengths
  - Less susceptible to noise and outliers
- Limitations
  - Biased towards globular clusters

## Hierarchical Clustering: Major Weaknesses

---

- Do not scale well ( $N$ : number of points)
  - Space complexity:
  - Time complexity:

## Hierarchical Clustering: Major Weaknesses

---

- Do not scale well ( $N$ : number of points)
  - Space complexity:  $O(N^2)$
  - Time complexity:  $O(N^3)$   
 $O(N^2 \log(N))$  for some cases/approaches
- Cannot undo what was done previously
- Quality varies in terms of distance measures
  - MIN (single link): susceptible to noise/outliers
  - MAX/GROUP AVERAGE: may not work well with non-globular clusters

# Cluster Analysis: Basic Concepts and Methods

---

- Cluster Analysis: Basic Concepts
- Similarity and distances
- Partitioning Methods
- Hierarchical Methods
- Density-Based Methods
- Probabilistic Methods
- Evaluation of Clustering

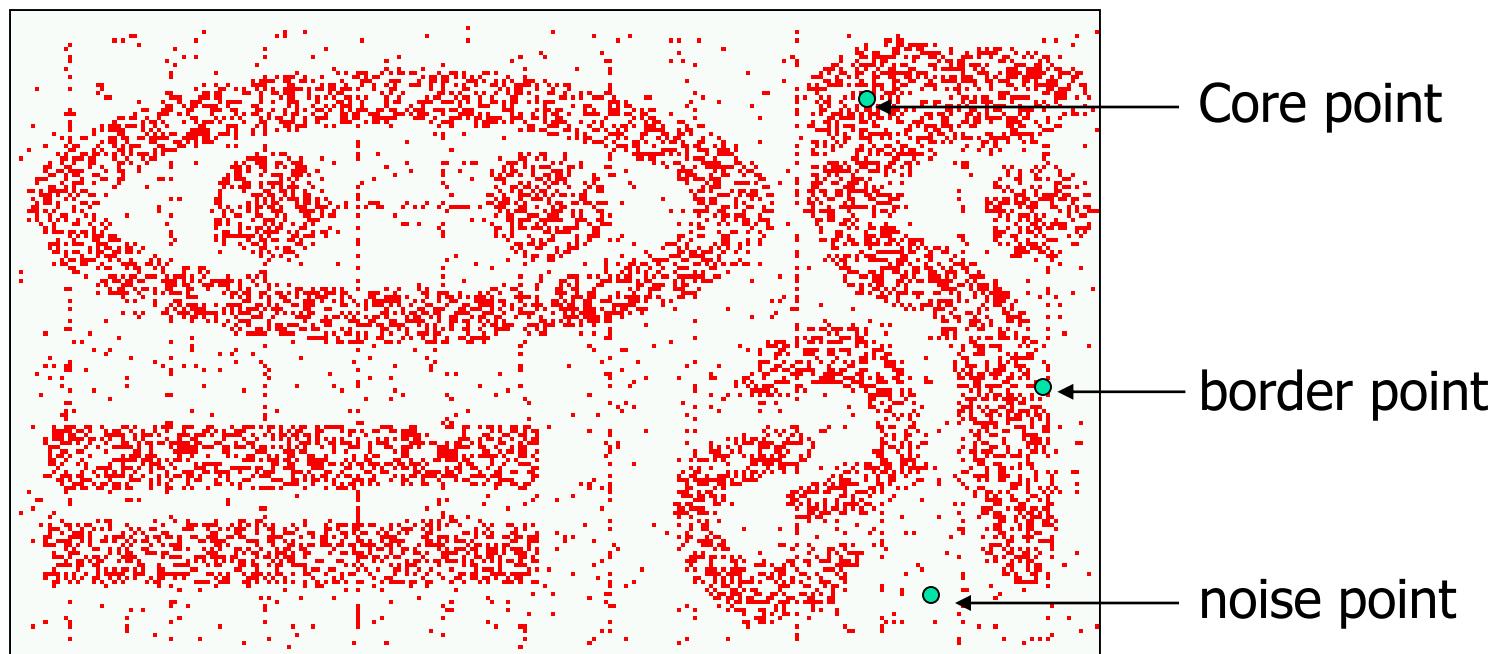
# Density-Based Clustering Methods

- Clustering based on density
- Major features:
  - Clusters of arbitrary shape
  - Handle noise
  - One scan
  - Need density parameters as termination condition
- Several interesting studies:
  - DBSCAN: Ester, et al. (KDD'96)
  - OPTICS: Ankerst, et al (SIGMOD'99).
  - DENCLUE: Hinneburg & D. Keim (KDD'98)
  - CLIQUE: Agrawal, et al. (SIGMOD'98) (more grid-based)



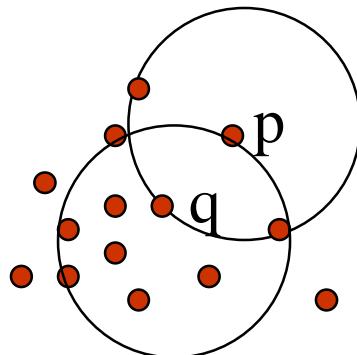
# DBSCAN: Basic Concepts

- Density = number of points within a specified radius
- **core point:** has high density
- **border point:** has less density, but in the neighborhood of a core point
- **noise point:** not a core point or a border point.



# DBScan: Definitions

- Two parameters:
  - *Eps*: radius of the neighbourhood
  - *MinPts*: Minimum number of points in an Eps-neighbourhood of that point
- $N_{Eps}(p)$ :  $\{q \text{ belongs to } D \mid \text{dist}(p,q) \leq Eps\}$
- core point:  $|N_{Eps}(q)| \geq MinPts$

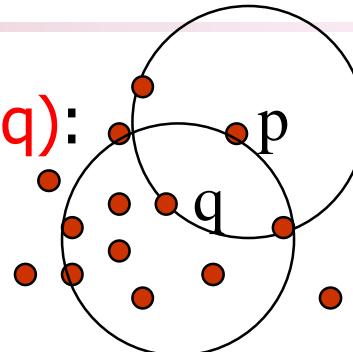


$MinPts = 5$

$Eps = 1 \text{ cm}$

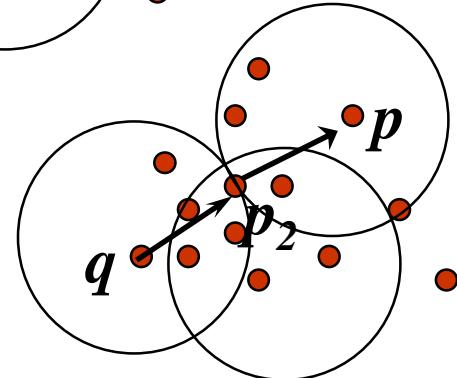
# DBScan: Definitions

- Directly density-reachable (p from q):  
 $p$  belongs to  $N_{Eps}(q)$

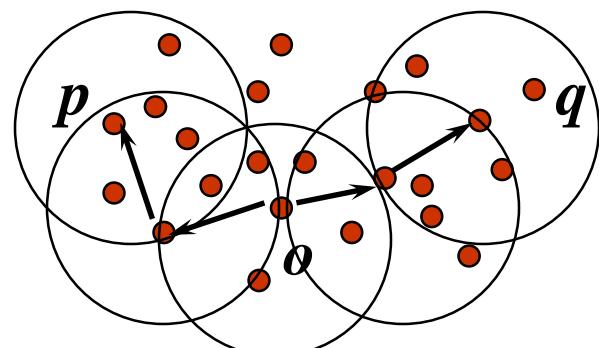


MinPts = 5  
Eps = 1 cm

- Density-reachable (p from q): if there is a chain of points  $p_1, \dots, p_n$ ,  $p_1 = q$ ,  $p_n = p$  such that  $p_{i+1}$  is directly density-reachable from  $p_i$

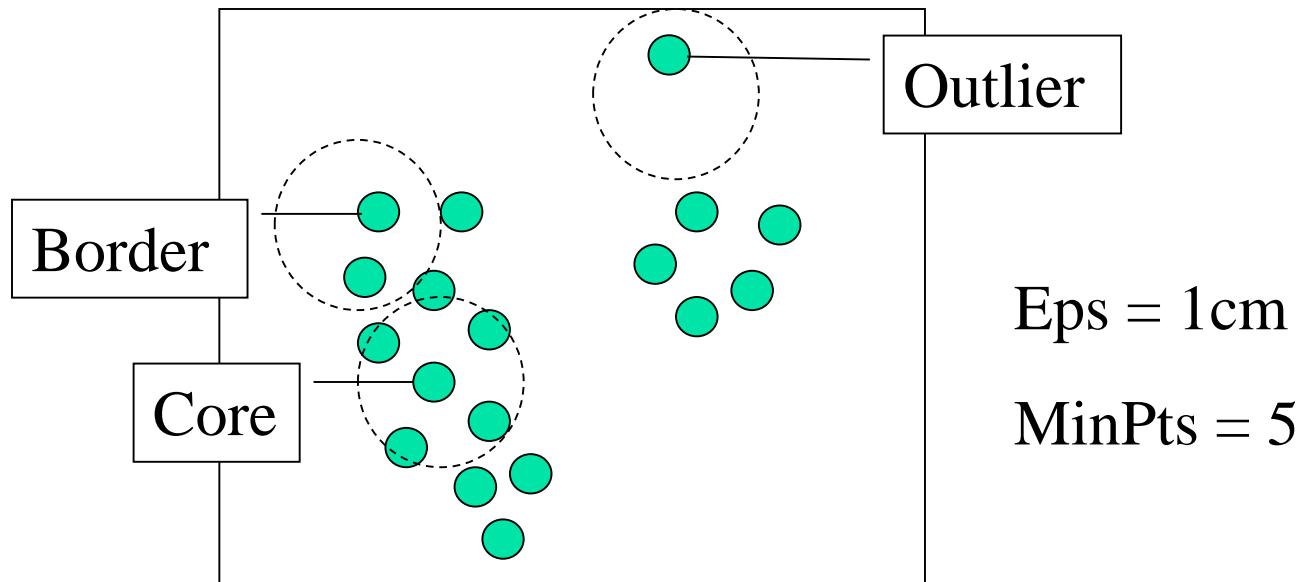


- Density-connected (p and q): if there is a point  $o$  such that both,  $p$  and  $q$  are density-reachable from  $o$  w.r.t.  
 $Eps$  and  $MinPts$



# DBSCAN: Cluster Definition

- A *cluster* is defined as a maximal set of density-connected points



# DBSCAN: The Algorithm

---

- Arbitrary select an unvisited point  $p$ , retrieve all neighbor points density-reachable from  $p$  w.r.t.  $Eps$  and  $MinPts$
- If  $p$  is a core point, a cluster is formed, add all neighbors of  $p$  to the cluster, and recursively add their neighbors if they are a core point
- Otherwise, mark  $p$  as a noise point
- Continue the process until all of the points have been processed.
- Complexity:  $O(n^2)$ . If a spatial index is used,  $O(n \log n)$

# DBSCAN: Sensitive to Parameters

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

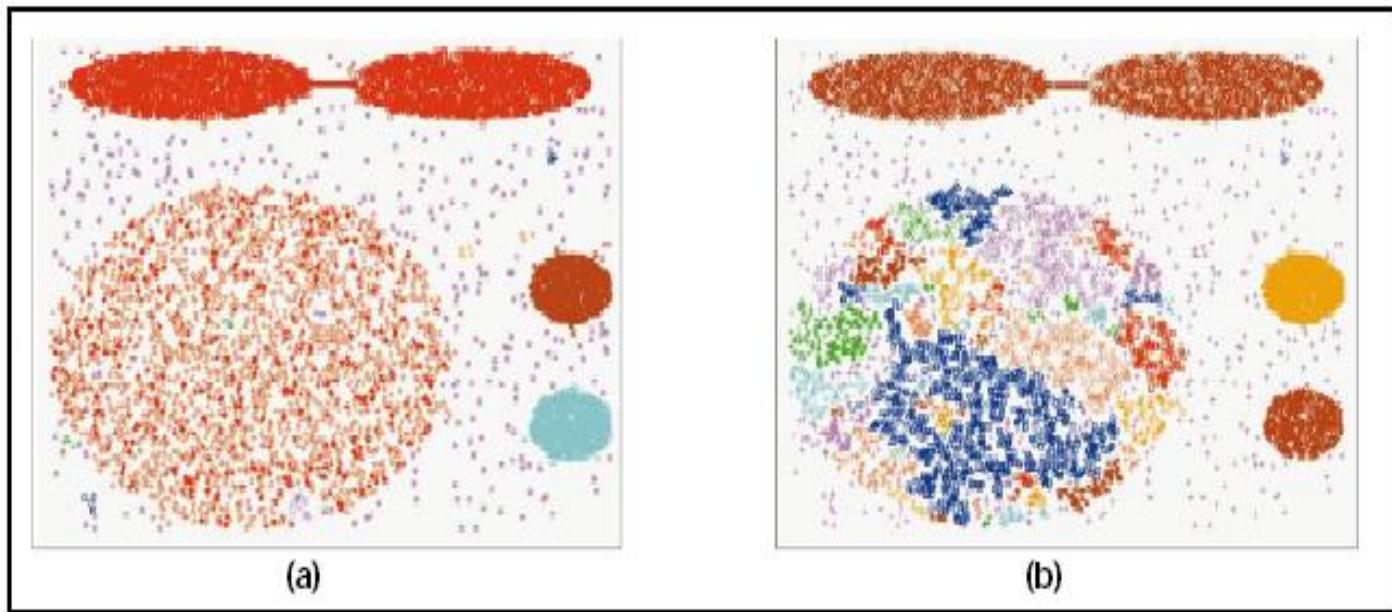
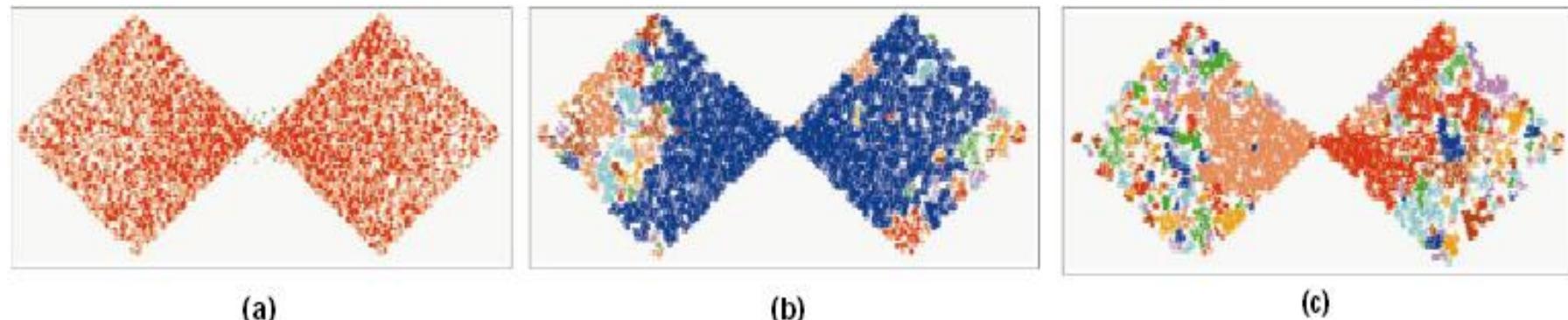
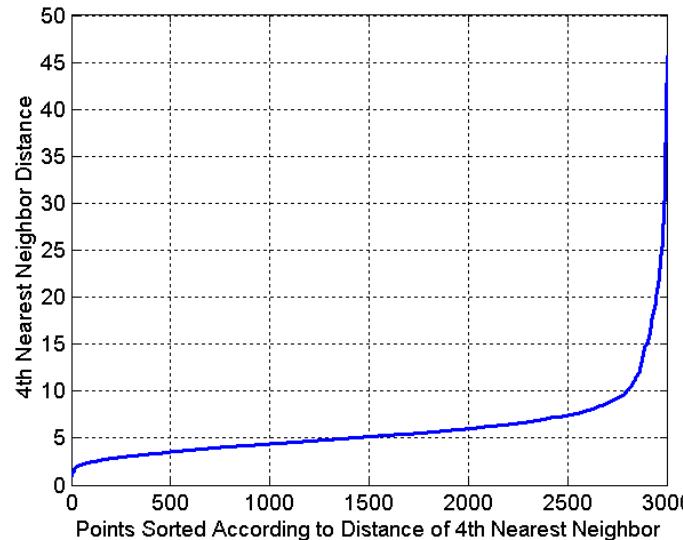


Figure 9. DBScan results for DS2 with MinPts at 4 and Eps at (a) 5.0, (b) 3.5, and (c) 3.0.



# DBSCAN: Determining EPS and MinPts

- Basic idea (given MinPts = k, find eps):
  - For points in a cluster, their  $k^{\text{th}}$  nearest neighbors are at roughly the same distance
  - Noise points have the  $k^{\text{th}}$  nearest neighbor at farther distance
- Plot sorted distance of every point to its  $k^{\text{th}}$  nearest neighbor



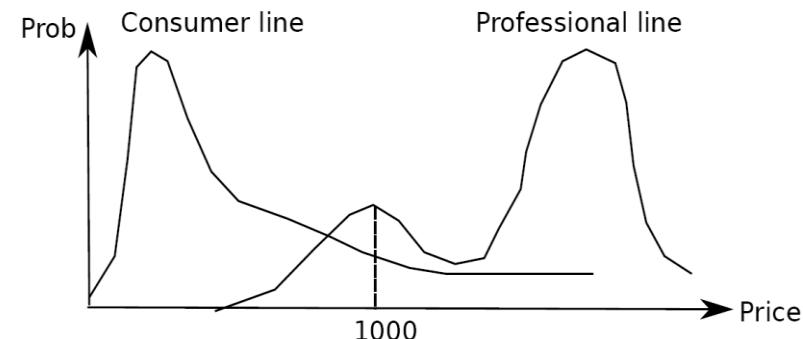
# Cluster Analysis: Basic Concepts and Methods

---

- Cluster Analysis: Basic Concepts
- Similarity and distances
- Partitioning Methods
- Hierarchical Methods
- Density-Based Methods
- Probabilistic Methods
- Evaluation of Clustering

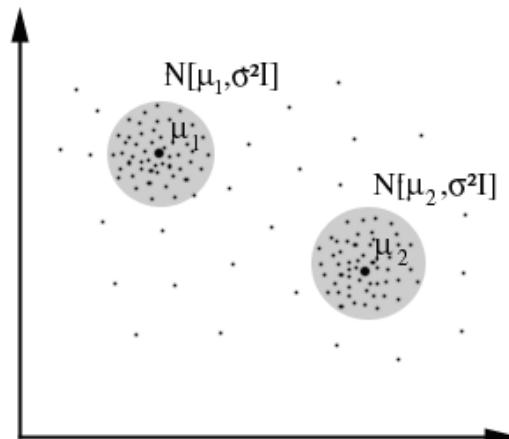
# Probabilistic Model-Based Clustering

- Data are instances of underlying hidden categories
- Cluster analysis is to find hidden categories.
- A hidden category is a distribution over the data space, which can be mathematically represented using a probability density function (or distribution function).
- Ex. 2 categories for digital cameras sold
  - consumer line vs. professional line
  - density functions  $f_1, f_2$  for  $C_1, C_2$
  - obtained by probabilistic clustering



# Clustering by Mixture Model

- A **mixture model** assumes data are generated by a mixture of probabilistic models
- Each cluster can be represented by a probabilistic model
  - e.g. a Gaussian (continuous) or a Poisson (discrete) distribution
- Data generation process: each observed object is generated independently
  - Choose a cluster,  $C_j$ , according to probabilities  $\omega_1, \dots, \omega_k$
  - Choose an instance of  $C_j$  according to its probability density function  $f_j$
- **Our task:** infer a set of  $k$  probabilistic models that is mostly likely to generate the data



# Model-Based Clustering

---

- A set  $C$  of  $k$  probabilistic clusters  $C_1, \dots, C_k$  with probability density functions  $f_1, \dots, f_k$ , respectively, and their probabilities  $\omega_1, \dots, \omega_k$ .
- Probability of an object  $o$  generated by cluster  $C_j$  is  $P(o|C_j) = \omega_j f_j(o)$
- Probability of  $o$  generated by the set of cluster  $C$  is  $P(o|C) = \sum_{j=1}^k \omega_j f_j(o)$
- Since objects are assumed to be generated independently, for a data set  $D = \{o_1, \dots, o_n\}$ , we have,

$$P(D|C) = \prod_{i=1}^n P(o_i|C) = \prod_{i=1}^n \sum_{j=1}^k \omega_j f_j(o_i)$$

- **Task:** Find a set  $C$  of  $k$  probabilistic clusters s.t.  $P(D|C)$  is maximized

# Univariate Gaussian Mixture Model

- $O = \{o_1, \dots, o_n\}$  ( $n$  observed objects),  $\Theta = \{\theta_1, \dots, \theta_k\}$  (parameters of the  $k$  distributions), and  $P_j(o_i | \theta_j)$  is the probability that  $o_i$  is generated from the  $j$ -th distribution using parameter  $\theta_j$ , we have

$$P(o_i | \Theta) = \sum_{j=1}^k \omega_j P_j(o_i | \Theta_j) \quad P(\mathbf{O} | \Theta) = \prod_{i=1}^n \sum_{j=1}^k \omega_j P_j(o_i | \Theta_j)$$

- Univariate Gaussian mixture model
  - Assume the probability density function of each cluster follows a 1-d Gaussian distribution. Suppose that there are  $k$  clusters with  $1/k$  prob.
  - The probability density function of each cluster are centered at  $\mu_j$  with standard deviation  $\sigma_j$ ,  $\theta_j = (\mu_j, \sigma_j)$ , we have

$$P(o_i | \Theta_j) = \frac{1}{\sqrt{2\pi}\sigma_j} e^{-\frac{(o_i - \mu_j)^2}{2\sigma_j^2}} \quad P(o_i | \Theta) = \sum_{j=1}^k \frac{1}{\sqrt{2\pi}\sigma_j} e^{-\frac{(o_i - \mu_j)^2}{2\sigma_j^2}}$$

$$P(\mathbf{O} | \Theta) = \prod_{i=1}^n \sum_{j=1}^k \frac{1}{\sqrt{2\pi}\sigma_j} e^{-\frac{(o_i - \mu_j)^2}{2\sigma_j^2}}$$

# The EM (Expectation Maximization) Algorithm

---

- **The (EM) algorithm:** A framework to approach maximum likelihood or maximum a posteriori estimates of parameters in statistical models.
  - **Expectation-step** assigns objects to clusters according to the current clustering or parameters of probabilistic clusters
  - **Maximization-step** finds the new clustering or parameters that maximize the expected likelihood

# Computing Mixture Models with EM

---

- Given  $n$  objects  $\mathbf{O} = \{o_1, \dots, o_n\}$ , we want to infer a set of parameters  $\Theta = \{\theta_1, \dots, \theta_k\}$  s.t.,  $P(\mathbf{O}|\Theta)$  is maximized, where  $\theta_j = (\mu_j, \sigma_j)$  are the mean and standard deviation of the  $j$ -th univariate Gaussian distribution
- We initially assign random values to parameters  $\theta_j$ , then iteratively conduct the Expectation (E) and Maximization (M) steps until converge
- At the **E-step**, for each object  $o_i$ , calculate the probability that  $o_i$  belongs to each distribution,

$$P(\Theta_j|o_i, \Theta) = \frac{P(o_i|\Theta_j)}{\sum_{l=1}^k P(o_i|\Theta_l)}$$

- At the **M-step**, adjust the parameters  $\theta_j = (\mu_j, \sigma_j)$  so that the expected likelihood  $P(\mathbf{O}|\Theta)$  is maximized

$$\mu_j = \sum_{i=1}^n o_i \frac{P(\Theta_j|o_i, \Theta)}{\sum_{l=1}^n P(\Theta_j|o_l, \Theta)} = \frac{\sum_{i=1}^n o_i P(\Theta_j|o_i, \Theta)}{\sum_{i=1}^n P(\Theta_j|o_i, \Theta)} \quad \sigma_j = \sqrt{\frac{\sum_{i=1}^n P(\Theta_j|o_i, \Theta)(o_i - \mu_j)^2}{\sum_{i=1}^n P(\Theta_j|o_i, \Theta)}}$$

# The EM (Expectation Maximization) Algorithm

---

- The k-means algorithm has two steps at each iteration:
  - **Expectation Step (E-step):** Given the current cluster centers, each object is assigned to the cluster whose center is closest to the object: An object is *expected to belong to the closest cluster*
  - **Maximization Step (M-step):** Given the cluster assignment, for each cluster, the algorithm *adjusts the center* so that *the sum of distance* from the objects assigned to this cluster and the new center is minimized

# Advantages and Disadvantages of Mixture Models

---

- Strength
  - Mixture models are more general than partitioning methods
  - Clusters can be characterized by a small number of parameters
  - The results may satisfy the statistical assumptions of the generative models
- Weakness
  - Converge to local optimal (overcome: run multi-times w. random initialization)
  - Computationally expensive if the number of distributions is large, or the data set contains very few observed data points
  - Need large data sets
  - Hard to estimate the number of clusters

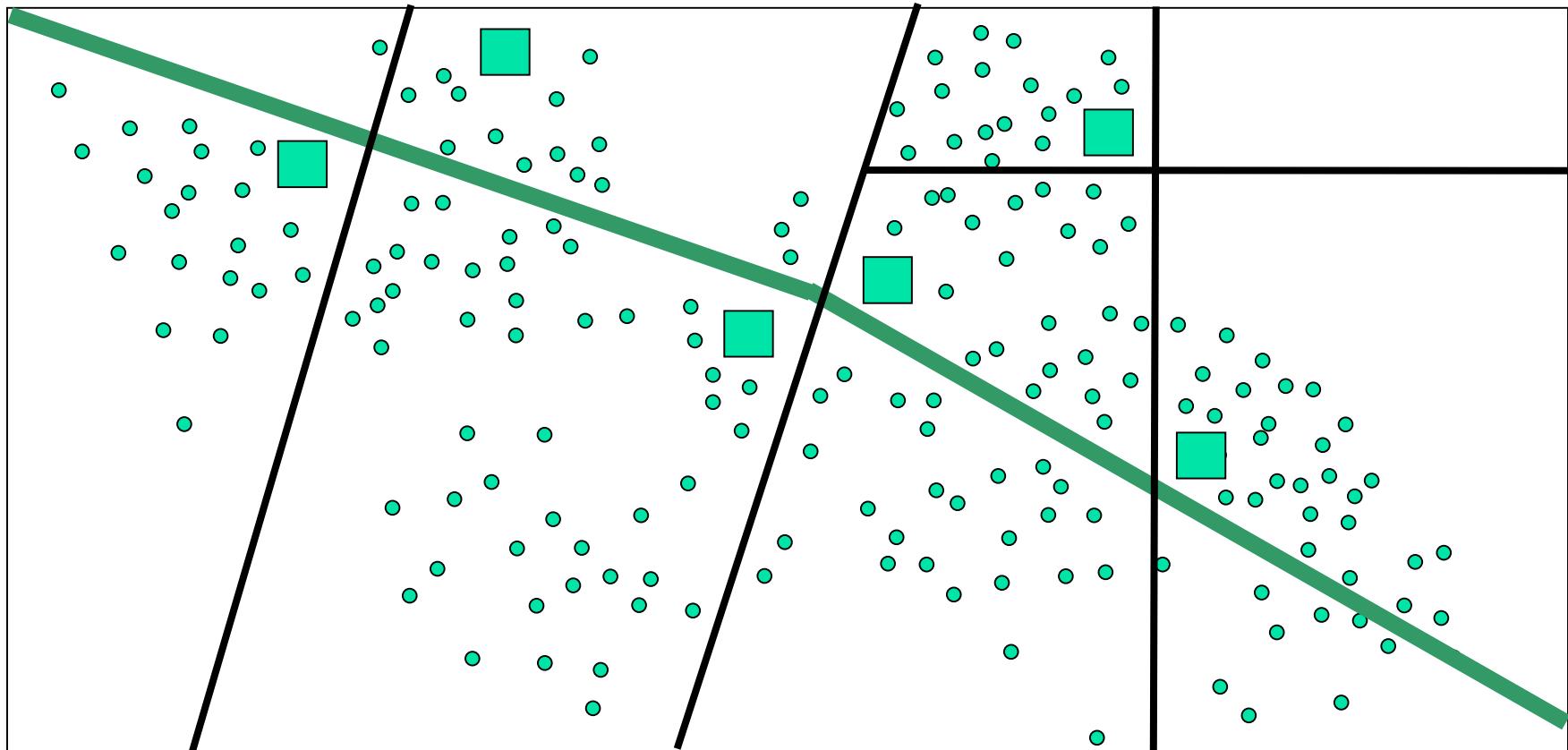
# Cluster Analysis: Basic Concepts and Methods

---

- Cluster Analysis: Basic Concepts
- Similarity and distances
- Partitioning Methods
- Hierarchical Methods
- Density-Based Methods
- Probabilistic Methods
- Evaluation of Clustering
- Clustering with constraints

# Why Constraint-Based Cluster Analysis?

- Need user feedback: Users know their applications the best
- Less parameters but more user-desired constraints, e.g., an ATM allocation problem: obstacle & desired clusters



# Categorization of Constraints

---

- Constraints on **instances**: specifies how a pair or a set of instances should be grouped in the cluster analysis
  - Must-link vs. cannot link constraints
    - $\text{must-link}(x, y)$ :  $x$  and  $y$  should be grouped into one cluster
  - Constraints can be defined using variables, e.g.,
    - $\text{cannot-link}(x, y)$  if  $\text{dist}(x, y) > d$
- Constraints on **clusters**: specifies a requirement on the clusters
  - E.g., specify the min # of objects in a cluster, the max diameter of a cluster, the shape of a cluster (e.g., a convex), # of clusters (e.g.,  $k$ )
- Constraints on **similarity** measurements: specifies a requirement that the similarity calculation must respect
  - E.g., driving on roads, obstacles (e.g., rivers, lakes)
- Issues: Hard vs. soft constraints; conflicting or redundant constraints

# Constraint-Based Clustering Methods (I): Handling Hard Constraints

---

- Handling hard constraints: Strictly respect the constraints in cluster assignments
- How to handle **must-link** and **cannot-link** constraints in k-means?

# Constraint-Based Clustering Methods (I): Handling Hard Constraints

---

- Handling hard constraints: Strictly respect the constraints in cluster assignments
- How to handle must-link and cannot-link constraints in k-means?
- Example: The COP-k-means algorithm
  - Generate super-instances for **must-link** objects
    - Compute the transitive closure of the must-link objects
    - Replace all objects in each subset by the mean
    - The super-instance also carries a weight, which is the number of objects it represents
  - Modified cluster assignment for **cannot-link** constraints
    - Modify the center-assignment process in k-means to a *nearest feasible center assignment*

# Constraint-Based Clustering Methods (II): Handling Soft Constraints

---

- Treated as an optimization problem:
  - When a clustering violates a soft constraint, a penalty is imposed on the clustering
- Overall objective:
  - Optimizing the clustering quality, and minimizing the constraint violation penalty
- Ex. CVQE (Constrained Vector Quantization Error) algorithm
  - Objective function: Sum of distance used in k-means, adjusted by the constraint violation penalties

# 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 be categorized into partitioning methods, hierarchical methods, density-based methods, and model-based methods