

Clustering Techniques (1)

Clustering Overview

Today's lecture

- What is clustering
- Partitional algorithms: K-means

Next lecture

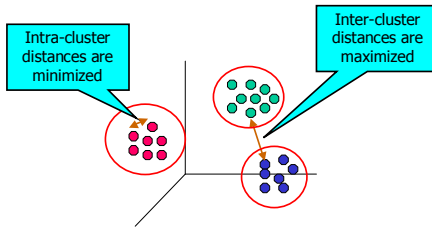
- Hierarchical algorithms
- Density-based algorithms: DBSCAN
- Techniques for clustering large databases
 - BIRCH
 - CURE

Data Mining: Clustering

2

What is Cluster Analysis?

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



Data Mining: Clustering

3

Applications of Cluster Analysis

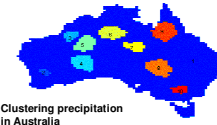
Understanding

- Group related documents for browsing, group genes and proteins that have similar functionality, or group stocks with similar price fluctuations

	Discovered Clusters	Industry Group
1	Applied Matl DOWN, Ray Network DOWN, COLLEGEN, Calsonic Sys DOWN, CSCO DOWN, J.P. DOWN, EMC Comm DOWN, INTEL DOWN, LSI Logic DOWN, Morgan Tech DOWN, Oracle Java DOWN, Tishco Inc DOWN, Natl Semiconduct DOWN, Veeva DOWN, SGE DOWN, Sun DOWN	Technology1-DOWN
2	Applied Genet DOWN, Amgen DOWN, BILLY DOWN, ADY Mater Device DOWN, Andros Corp DOWN, Computer Assoc DOWN, Genentx Corp DOWN, Compag DOWN, EMC Corp DOWN, Gen Inc DOWN, Minerals DOWN, Microsoft DOWN, Scientific DOWN	Technology2-DOWN
3	Farmvix Mac DOWN, Fiat Home Loan DOWN, MBIA Corp DOWN, Morgan Stanley DOWN	Financial-DOWN
4	Baker Hughes UP, Denver Steel UP, Halliburton RE UP, Lennovo Lead UP, Phillips Petrol UP, Unocal UP, Schlumberger UP	OIL-UP

Summarization

- Reduce the size of large data sets



Clustering precipitation in Australia

Data Mining: Clustering

4

Clustering Examples and Applications

- Segment customer database based on similar buying patterns
- Group houses in a town into neighborhoods based on similar features
- Group plants into categories
- Identify similar Web usage patterns

Data Mining: Clustering

5

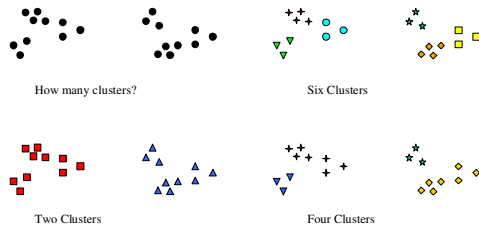
What is not Cluster Analysis?

- Supervised classification
 - Have class label information
- Simple segmentation
 - Dividing students into different registration groups alphabetically, by last name
- Group results of a query
 - When groupings are a result of an external specification
- Graph partitioning

Data Mining: Clustering

6

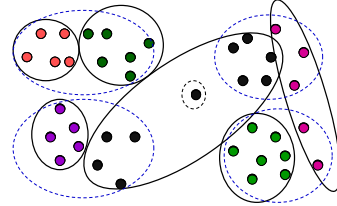
Notion of a Cluster can be Ambiguous



Data Mining: Clustering

7

Clustering can be performed in various ways



Attribute Value Based Clustering
Geographic Distance Based Clustering

Data Mining: Clustering

8

Clustering Problem

- Given a database $D = \{t_1, t_2, \dots, t_n\}$ of tuples and possibly also an integer value k , the **clustering problem** is to define a mapping $f: D \rightarrow \{1, \dots, k\}$ such that each t_i is assigned to one cluster K_j , $1 \leq j \leq k$.
- A **cluster**, K_j , contains precisely those tuples mapped to it.

Data Mining: Clustering

9

Clustering vs. Classification

- In **classification** the set of groups to which objects belong is given, and the task is to discriminate between groups on the basis of values for the attributes of the objects
- In **clustering** there is no prior knowledge: the meaning (and often the number) of clusters is unknown and the objective is to discover them
- Clustering is sometimes known as *unsupervised learning*

Data Mining: Clustering

10

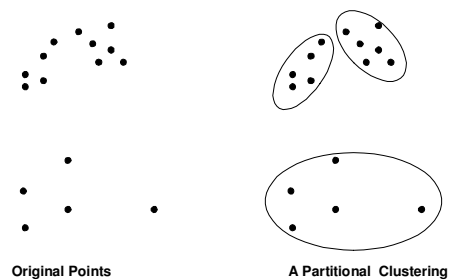
Types of Clusterings

- 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

Data Mining: Clustering

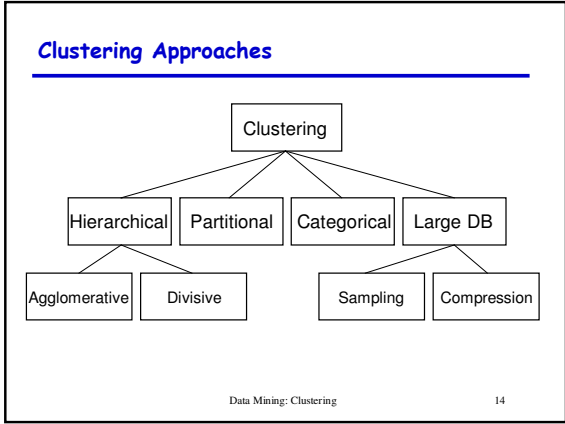
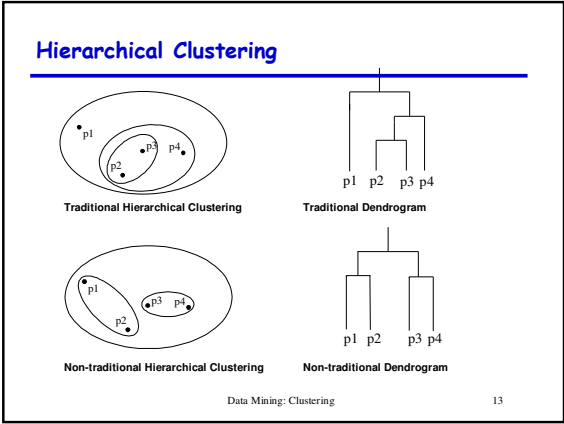
11

Partitional Clustering

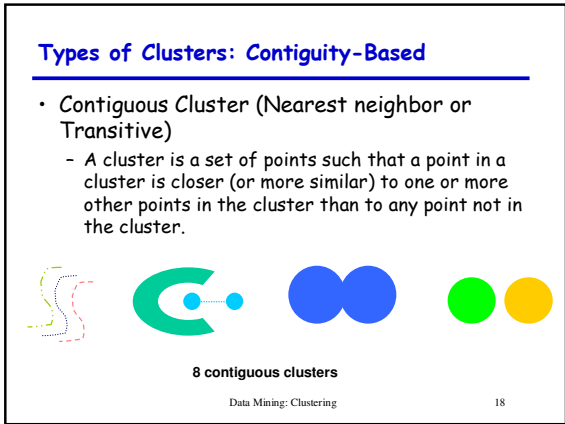
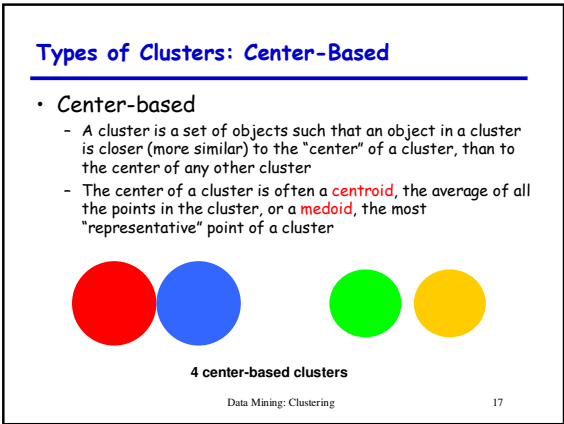
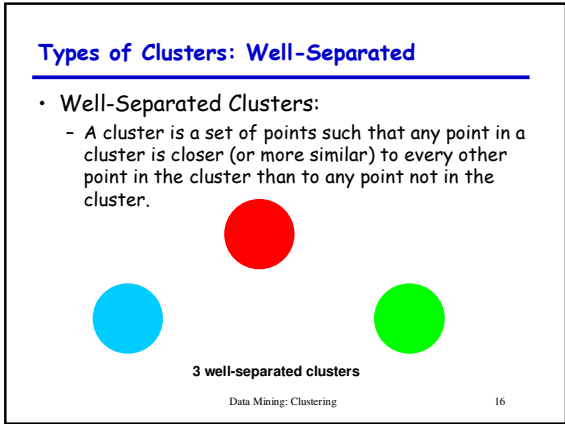


Data Mining: Clustering

12

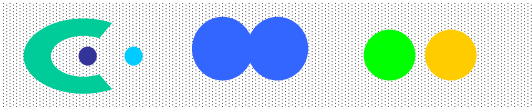


- ### Other Distinctions Between Sets of Clusters
- **Exclusive vs. non-exclusive**
 - In non-exclusive clusterings, points may belong to multiple clusters.
 - Can represent multiple classes or 'border' points
 - **Fuzzy vs. non-fuzzy**
 - In fuzzy clustering, a point belongs to every cluster with some weight between 0 and 1
 - Weights must sum to 1
 - Probabilistic clustering has similar characteristics
 - **Partial vs. complete**
 - In some cases, we only want to cluster some of the data
 - **Heterogeneous vs. homogeneous**
 - Cluster of widely different sizes, shapes, and densities
- Data Mining: Clustering 15



Types of Clusters: Density-Based

- Density-based
 - A cluster is a dense region of points, which is separated by low-density regions, from other regions of high density.
 - Used when the clusters are irregular or intertwined, and when noise and outliers are present.



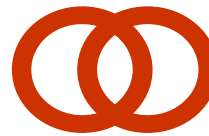
6 density-based clusters

Data Mining: Clustering

19

Types of Clusters: Conceptual Clusters

- Shared Property or Conceptual Clusters
 - Finds clusters that share some common property or represent a particular concept.



2 Overlapping Circles

Data Mining: Clustering

20

Types of Clusters: Objective Function

- Clusters Defined by an Objective Function
 - Finds clusters that minimize or maximize an objective function.
 - Enumerate all possible ways of dividing the points into clusters and evaluate the 'goodness' of each potential set of clusters by using the given objective function. (NP Hard)
 - Can have global or local objectives.
 - Hierarchical clustering algorithms typically have local objectives
 - Partitional algorithms typically have global objectives
 - A variation of the global objective function approach is to fit the data to a parameterized model.
 - Parameters for the model are determined from the data.
 - Mixture models assume that the data is a 'mixture' of a number of statistical distributions.

Data Mining: Clustering

21

Types of Clusters: Objective Function (Cont)

- Map the clustering problem to a different domain and solve a related problem in that domain
 - Proximity matrix defines a weighted graph, where the nodes are the points being clustered, and the weighted edges represent the proximities between points.
 - Clustering is equivalent to breaking the graph into connected components, one for each cluster.
 - Want to minimize the edge weight between clusters and maximize the edge weight within clusters.

Data Mining: Clustering

22

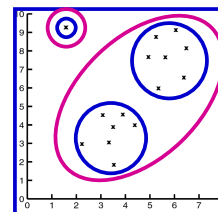
Characteristics of the Data are Important

- Type of proximity or density measure
 - This is a derived measure, but central to clustering
- Sparseness
 - Dictates type of similarity
 - Adds to efficiency
- Attribute type
 - Dictates type of similarity
- Type of Data
 - Dictates type of similarity
 - Other characteristics, e.g., autocorrelation
- Dimensionality
- Noise and Outliers
- Type of Distribution

Data Mining: Clustering

23

Impact of Outliers on Clustering



Data Mining: Clustering

24

Cluster Parameters

$$centroid = C_m = \frac{\sum_{i=1}^N (t_{mi})}{N}$$

$$radius = R_m = \sqrt{\frac{\sum_{i=1}^N (t_{mi} - C_m)^2}{N}}$$

$$diameter = D_m = \sqrt{\frac{\sum_{i=1}^N \sum_{j=1}^N (t_{mi} - t_{mj})^2}{(N)(N-1)}}$$

Data Mining: Clustering

25

Clustering Algorithms

- Partitional algorithms
 - K-means and its variants
- Hierarchical clustering algorithms
- Density-based clustering algorithms

Data Mining: Clustering

26

K-means Clustering

- Partitional clustering approach
- Each cluster is associated with a **centroid** (center point)
- Each point is assigned to the cluster with the closest centroid
- Number of clusters, K, must be specified
- The basic algorithm is very simple

- 1: Select K points as the initial centroids.
- 2: **repeat**
- 3: Form K clusters by assigning all points to the closest centroid.
- 4: Recompute the centroid of each cluster.
- 5: **until** The centroids don't change

Data Mining: Clustering

27

K-means Clustering - Example

- Given a cluster $K_i = \{t_{i1}, t_{i2}, \dots, t_{im}\}$, let the **cluster mean** be $m_i = (1/m)(t_{i1} + \dots + t_{im})$

Given: {2,4,10,12,3,20,30,11,25}, k=2

- Randomly pick some initial means: $m_1=3, m_2=4$
- $K_1=\{2,3\}, K_2=\{4,10,12,20,30,11,25\}, m_1=2.5, m_2=16$
- $K_1=\{2,3,4\}, K_2=\{10,12,20,30,11,25\}, m_1=3, m_2=18$
- $K_1=\{2,3,4,10\}, K_2=\{12,20,30,11,25\}, m_1=4.75, m_2=19.6$
- $K_1=\{2,3,4,10,11,12\}, K_2=\{20,30,25\}, m_1=7, m_2=25$

Stop as the clusters with these means are the same.

Data Mining: Clustering

28

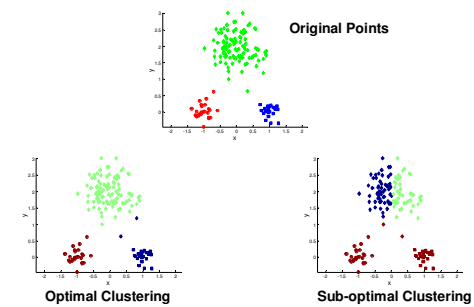
K-means Clustering - Details

- 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.
- Convergence happens in the first few iterations.
 - Often the stopping condition is changed to 'until relatively few points change clusters'
- Complexity is $O(n * K * I * d)$
 - n = number of points, K = number of clusters,
 - I = number of iterations, d = number of attributes

Data Mining: Clustering

29

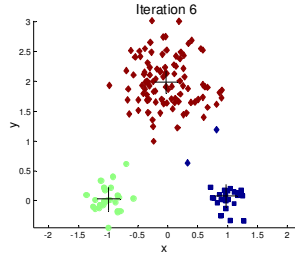
Two different K-means Clusterings



Data Mining: Clustering

30

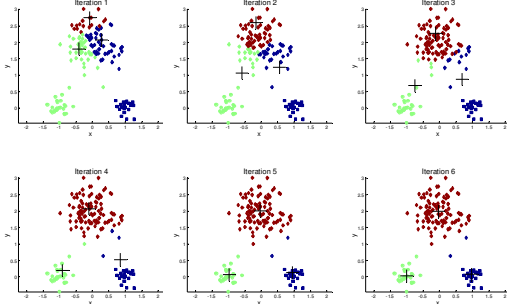
Importance of Choosing Initial Centroids



Data Mining: Clustering

31

Importance of Choosing Initial Centroids



Data Mining: Clustering

32

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.

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

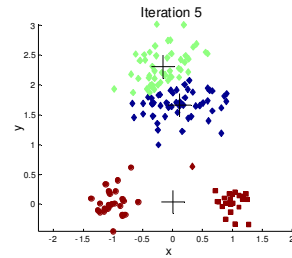
- x is a data point in cluster C_i and m_i is the representative point for cluster C_i
 - can show that m_i corresponds to the center (mean) of the cluster
- Given two clusters, we can choose the one with the smallest error
- One easy way to reduce SSE is to increase K , the number of clusters

A good clustering with smaller K can have a lower SSE than a poor clustering with higher K

Data Mining: Clustering

33

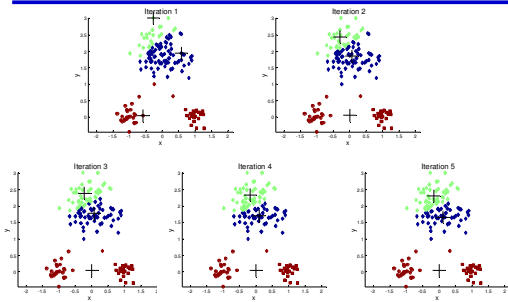
Importance of Choosing Initial Centroids



Data Mining: Clustering

34

Importance of Choosing Initial Centroids



Data Mining: Clustering

35

Problems with Selecting Initial Points

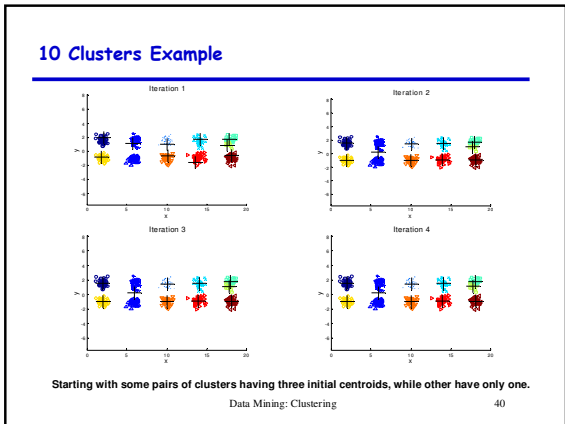
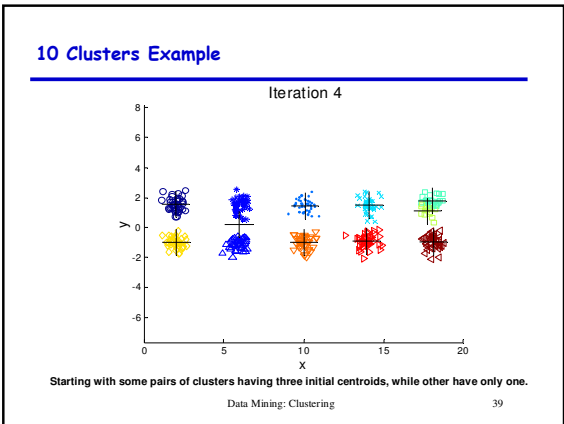
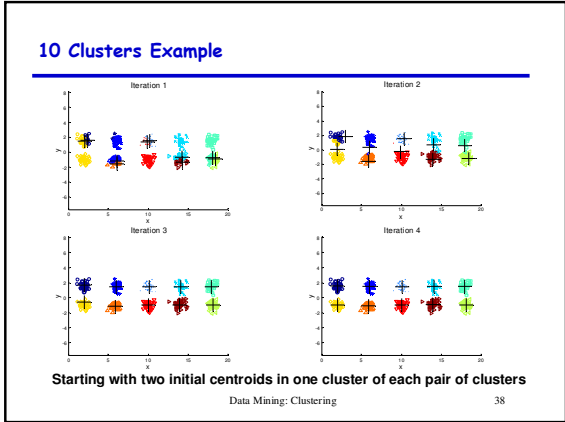
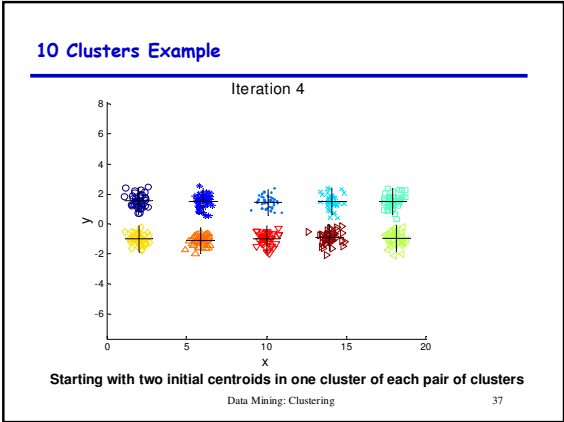
- If there are K 'real' clusters then the chance of selecting one centroid from each cluster is small.
 - Chance is relatively small when K is large
 - If clusters are the same size, n , then

$$P = \frac{\text{number of ways to select one centroid from each cluster}}{\text{number of ways to select } K \text{ centroids}} = \frac{K!n^K}{(Kn)^K} = \frac{K!}{K^K}$$

- For example, if $K = 10$, then probability = $10!/10^{10} = 0.00036$
- Sometimes the initial centroids will readjust themselves in 'right' way, and sometimes they don't
- Consider an example of five pairs of clusters

Data Mining: Clustering

36



- ### 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
- Data Mining: Clustering 41

- ### 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
 - Can use these steps during the clustering process
- Data Mining: Clustering 42

Bisecting K-means

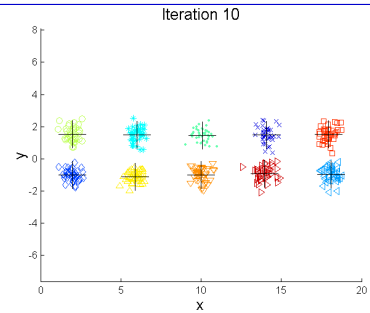
- Bisecting K-means algorithm
 - Variant of K-means that can produce a partitional or a hierarchical clustering

```
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
```

Data Mining: Clustering

43

Bisecting K-means Example



Data Mining: Clustering

44

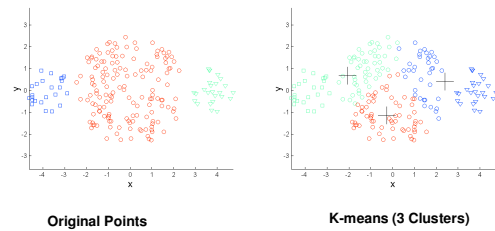
Limitations of K-means

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

Data Mining: Clustering

45

Limitations of K-means: Differing Sizes



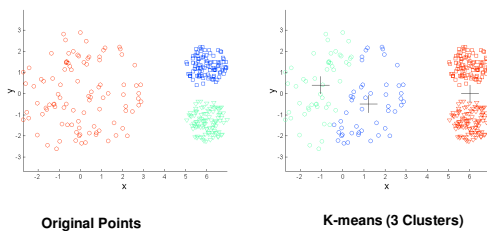
Original Points

K-means (3 Clusters)

Data Mining: Clustering

46

Limitations of K-means: Differing Density



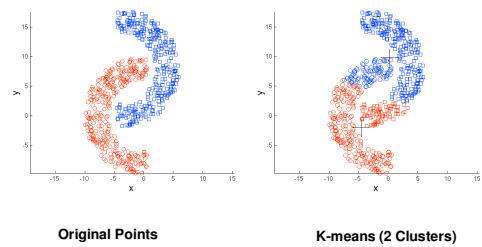
Original Points

K-means (3 Clusters)

Data Mining: Clustering

47

Limitations of K-means: Non-globular Shapes



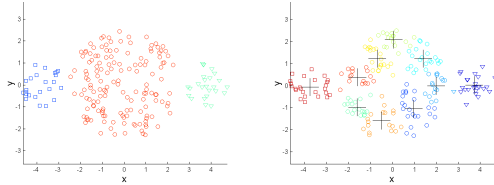
Original Points

K-means (2 Clusters)

Data Mining: Clustering

48

Overcoming K-means Limitations

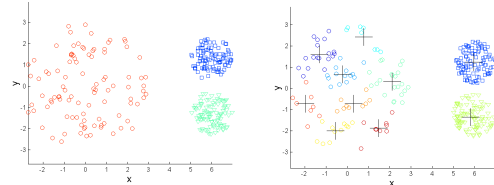


Original Points

K-means Clusters

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

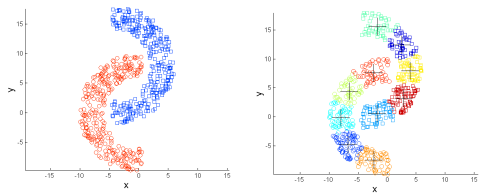
Overcoming K-means Limitations



Original Points

K-means Clusters

Overcoming K-means Limitations



Original Points

K-means Clusters