



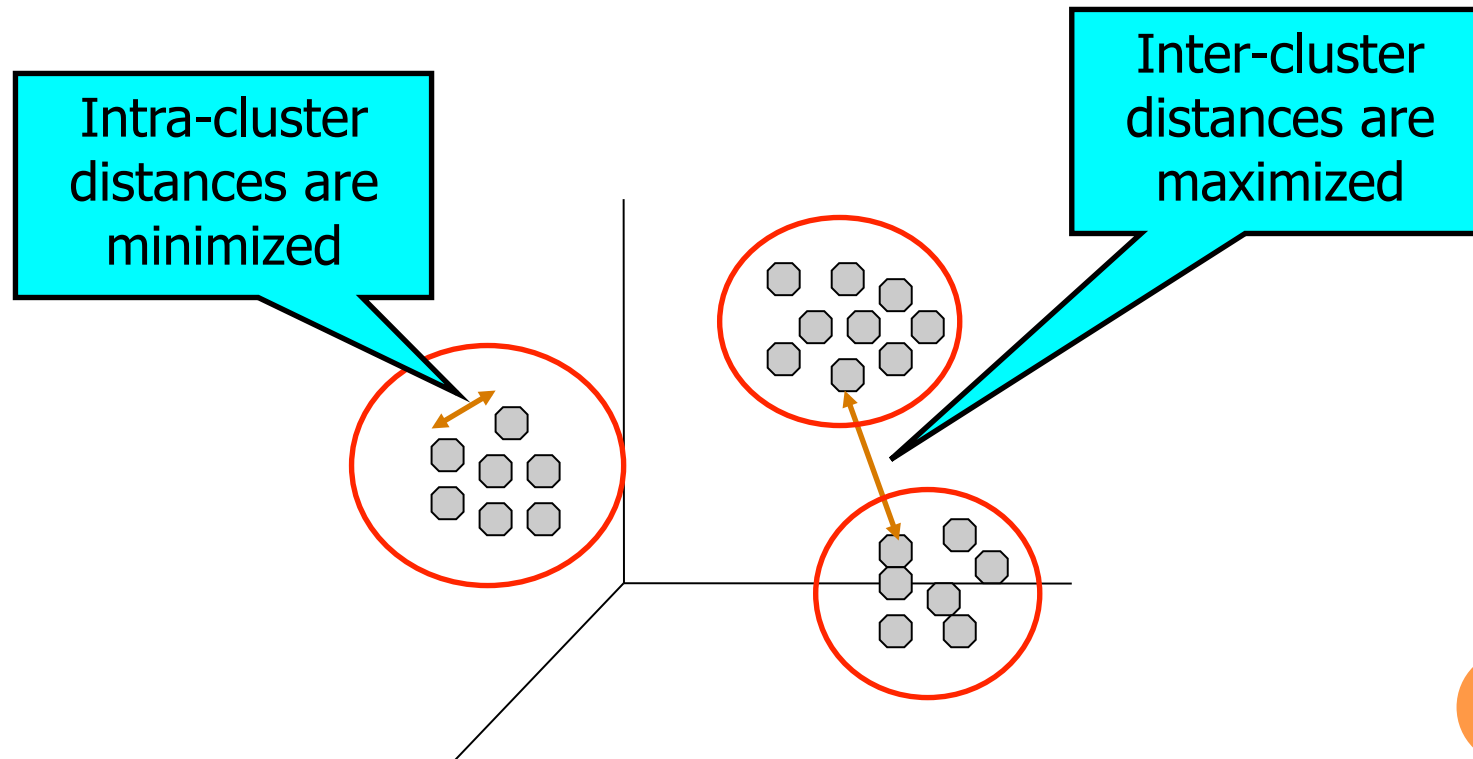
DATA MINING

CLUSTER ANALYSIS: BASIC CONCEPTS AND ALGORITHMS

Chiara Renso
KDD-LAB
ISTI- CNR, Pisa, Italy

WHAT IS 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



CLUSTERING: APPLICATION 1

- Market Segmentation:
 - Goal: subdivide a market into distinct subsets of customers where any subset may conceivably be selected as a market target to be reached with a distinct marketing mix.
 - Approach:
 - Collect different attributes of customers based on their geographical and lifestyle related information.
 - Find clusters of similar customers.
 - Measure the clustering quality by observing buying patterns of customers in same cluster vs. those from different clusters.

CLUSTERING: APPLICATION 2

- Document Clustering:
 - Goal: To find groups of documents that are similar to each other based on the important terms appearing in them.
 - Approach: To identify frequently occurring terms in each document. Form a similarity measure based on the frequencies of different terms. Use it to cluster.
 - Gain: Information Retrieval can utilize the clusters to relate a new document or search term to clustered documents.

ILLUSTRATING DOCUMENT CLUSTERING

- Clustering Points: 3204 Articles of Los Angeles Times.
- Similarity Measure: How many words are common in these documents (after some word filtering).

<i>Category</i>	<i>Total Articles</i>	<i>Correctly Placed</i>
<i>Financial</i>	555	364
<i>Foreign</i>	341	260
<i>National</i>	273	36
<i>Metro</i>	943	746
<i>Sports</i>	738	573
<i>Entertainment</i>	354	278

SIMILARITY

SIMILARITY AND DISSIMILARITY

- Similarity
 - Numerical measure of how alike two data objects are.
 - Is higher when objects are more alike.
 - Often falls in the range $[0,1]$
- Dissimilarity
 - Numerical measure of how different are two data objects
 - Lower when objects are more alike
 - Minimum dissimilarity is often 0
 - Upper limit varies
- Proximity refers to a similarity or dissimilarity

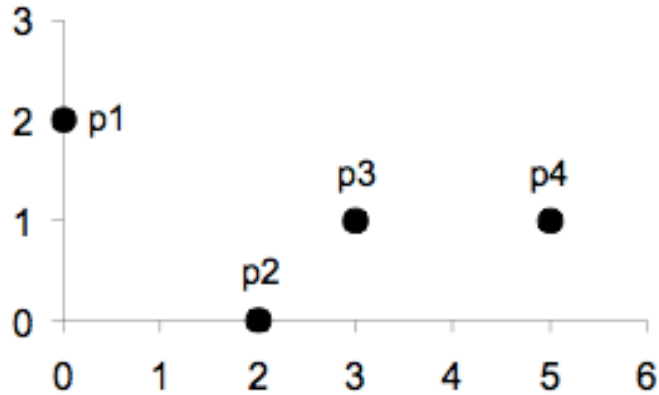
EUCLIDEAN DISTANCE

- Euclidean Distance

$$dist = \sqrt{\sum_{k=1}^n (p_k - q_k)^2}$$

- Where n is the number of dimensions (attributes) and p_k and q_k are, respectively, the k^{th} attributes (components) or data objects p and q .
- Standardization is necessary, if scales differ.

EUCLIDEAN DISTANCE



point	x	y
p1	0	2
p2	2	0
p3	3	1
p4	5	1

	p1	p2	p3	p4
p1	0	2.828	3.162	5.099
p2	2.828	0	1.414	3.162
p3	3.162	1.414	0	2
p4	5.099	3.162	2	0

Distance Matrix

MINKOWSKI DISTANCE

- Minkowski Distance is a generalization of Euclidean Distance

$$dist = \left(\sum_{k=1}^n |p_k - q_k|^r \right)^{\frac{1}{r}}$$

Where r is a parameter, n is the number of dimensions (attributes) and p_k and q_k are, respectively, the k th attributes (components) or data objects p and q .

MINKOWSKI DISTANCE: EXAMPLES

- $r = 1$. City block (Manhattan, taxicab, L_1 norm) distance.
 - A common example of this is the Hamming distance, which is just the number of bits that are different between two binary vectors
- $r = 2$. Euclidean distance
- $r \rightarrow \infty$. “supremum” (L_{\max} norm, L_∞ norm) distance.
 - This is the maximum difference between any component of the vectors
- Do not confuse r with n , i.e., all these distances are defined for all numbers of dimensions.

COMMON PROPERTIES OF A DISTANCE

- Distances, such as the Euclidean distance, have some well known properties.

1. $d(p, q) \geq 0$ for all p and q and $d(p, q) = 0$ only if $p = q$. (Positive definiteness)
2. $d(p, q) = d(q, p)$ for all p and q . (Symmetry)
3. $d(p, r) \leq d(p, q) + d(q, r)$ for all points p, q , and r . (Triangle Inequality)

where $d(p, q)$ is the distance (dissimilarity) between points (data objects), p and q .

- A distance that satisfies these properties is a **metric**

COMMON PROPERTIES OF A SIMILARITY

- Similarities, also have some well known properties.
 1. $s(p, q) = 1$ (or maximum similarity) only if $p = q$.
 2. $s(p, q) = s(q, p)$ for all p and q . (Symmetry)

where $s(p, q)$ is the similarity between points (data objects), p and q .

SIMILARITY BETWEEN BINARY VECTORS

- Common situation is that objects, p and q , have only binary attributes
- Compute similarities using the following quantities

M_{01} = the number of attributes where p was 0 and q was 1

M_{10} = the number of attributes where p was 1 and q was 0

M_{00} = the number of attributes where p was 0 and q was 0

M_{11} = the number of attributes where p was 1 and q was 1

- Simple Matching and Jaccard Coefficients

SMC = number of matches / number of attributes

$$= (M_{11} + M_{00}) / (M_{01} + M_{10} + M_{11} + M_{00})$$

J = number of 11 matches / number of not-both-zero attributes values

$$= (M_{11}) / (M_{01} + M_{10} + M_{11})$$

SMC VERSUS JACCARD: EXAMPLE

- $p = 1\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0$

- $q = 0\ 0\ 0\ 0\ 0\ 0\ 1\ 0\ 0\ 1$

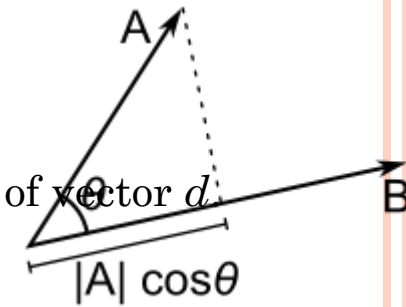
- $M_{01} = 2$ (the number of attributes where p was 0 and q was 1)
- $M_{10} = 1$ (the number of attributes where p was 1 and q was 0)
- $M_{00} = 7$ (the number of attributes where p was 0 and q was 0)
- $M_{11} = 0$ (the number of attributes where p was 1 and q was 1)

COSINE SIMILARITY

- If d_1 and d_2 are two document vectors, then

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

where \cdot indicates vector dot product and $||d||$ is the length of vector d .



- Example:

$$d_1 = 3 \ 2 \ 0 \ 5 \ 0 \ 0 \ 0 \ 2 \ 0 \ 0$$

$$d_2 = 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 2$$

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

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

$$||d_2|| = (1^2 + 0^2 + 0^2 + 0^2 + 0^2 + 0^2 + 0^2 + 1^2 + 0^2 + 2^2)^{0.5} = (6)^{0.5} = 2.245$$

$$\cos(d_1, d_2) = .3150$$

$$A \cdot B = ||A|| ||B|| \cos \theta$$

Pitagora Theorem

CORRELATION

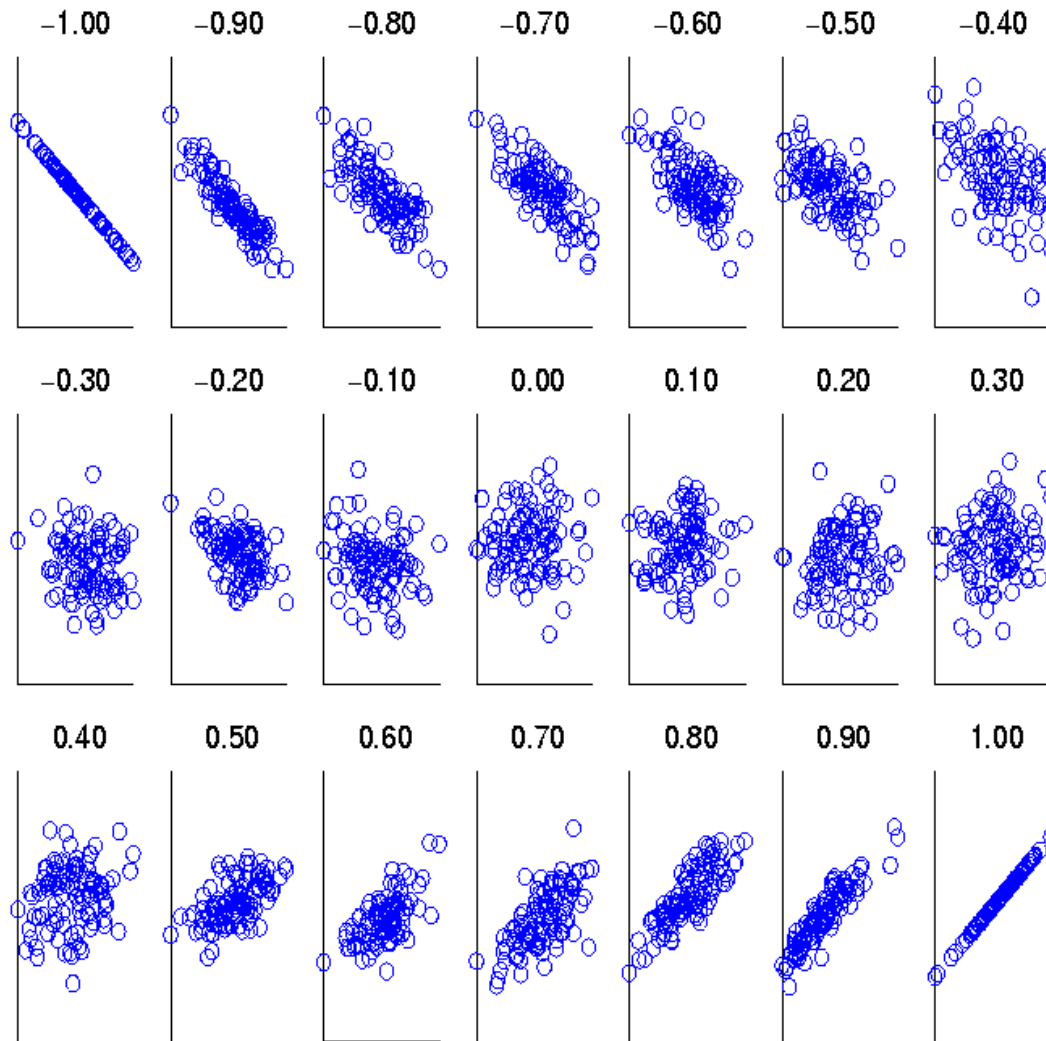
- Correlation measures the linear relationship between objects
- To compute correlation, we standardize data objects, p and q , and then take their dot product

$$p'_k = (p_k - \text{mean}(p)) / \text{std}(p)$$

$$q'_k = (q_k - \text{mean}(q)) / \text{std}(q)$$

$$\text{correlation}(p, q) = p' \cdot q'$$

VISUALLY EVALUATING CORRELATION



Scatter plots showing the similarity from -1 to 1 .

DENSITY

- Density-based clustering require a notion of density
- Examples:
 - Euclidean density
 - ◆ Euclidean density = number of points per unit volume
 - Probability density
 - Graph-based density

EUCLIDEAN DENSITY – CENTER-BASED

- Euclidean density is the number of points within a specified radius of the point

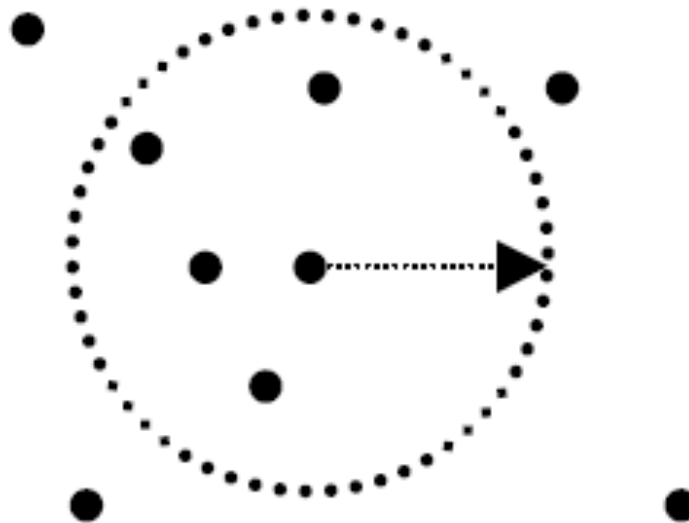


Figure 7.14. Illustration of center-based density.

CLUSTERING TECHNIQUES

21

APPLICATIONS OF CLUSTER ANALYSIS

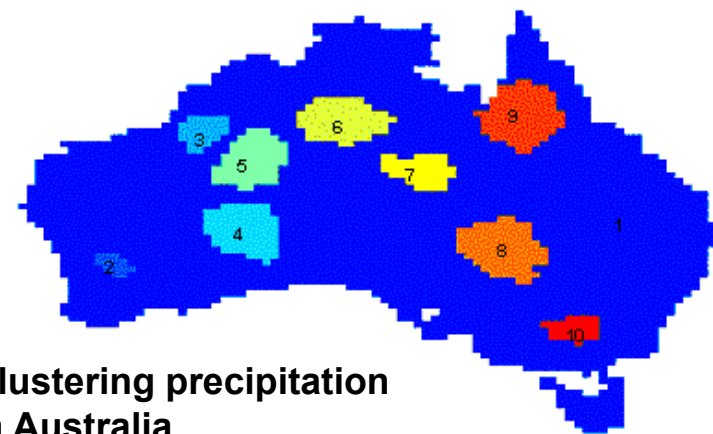
● Understanding

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

	<i>Discovered Clusters</i>	<i>Industry Group</i>
1	Applied-Matl-DOWN, Bay-Network-Down, 3-COM-DOWN, Cabletron-Sys-DOWN, CISCO-DOWN, HP-DOWN, DSC-Comm-DOWN, INTEL-DOWN, LSI-Logic-DOWN, Micron-Tech-DOWN, Texas-Inst-Down, Tellabs-Inc-Down, Natl-Semiconduct-DOWN, Oracle-DOWN, SGI-DOWN, Sun-DOWN	Technology 1-DOWN
2	Apple-Cmp-DOWN, Autodesk-DOWN, DEC-DOWN, ADV-Micro-Device-DOWN, Andrew-Corp-DOWN, Computer-Assoc-DOWN, Circuit-City-DOWN, Compaq-DOWN, EMC-Corp-DOWN, Gen-Inst-DOWN, Motorola-DOWN, Microsoft-DOWN, Scientific-Atl-DOWN	Technology 2-DOWN
3	Fannie-Mae-DOWN, Fed-Home-Loan-DOWN, MBNA-Corp-DOWN, Morgan-Stanley-DOWN	Financial-DOWN
4	Baker-Hughes-UP, Dresser-Inds-UP, Halliburton-HLD-UP, Louisiana-Land-UP, Phillips-Petro-UP, Unocal-UP, Schlumberger-UP	Oil-UP

● Summarization

- Reduce the size of large data sets

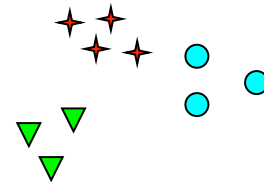


Clustering precipitation in Australia

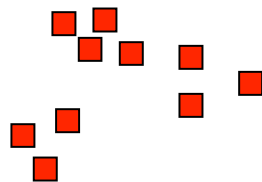
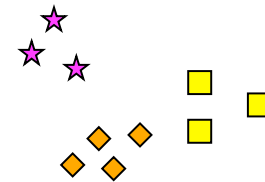
NOTION OF A CLUSTER CAN BE AMBIGUOUS



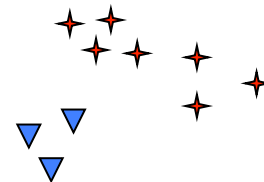
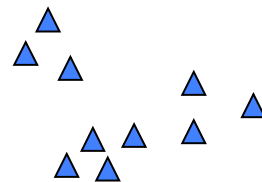
How many clusters?



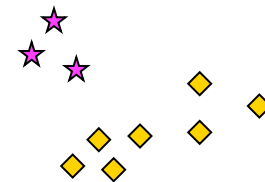
Six Clusters



Two Clusters



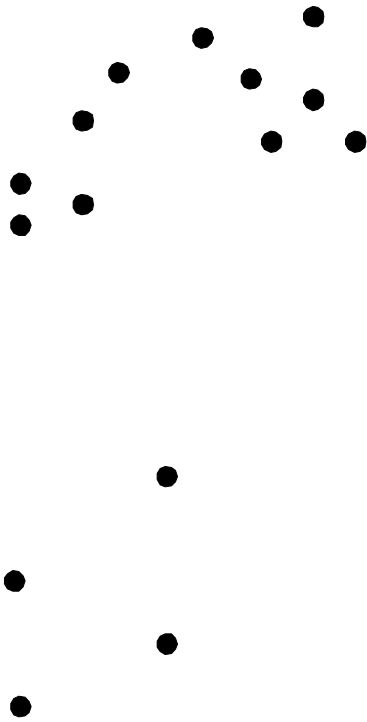
Four Clusters



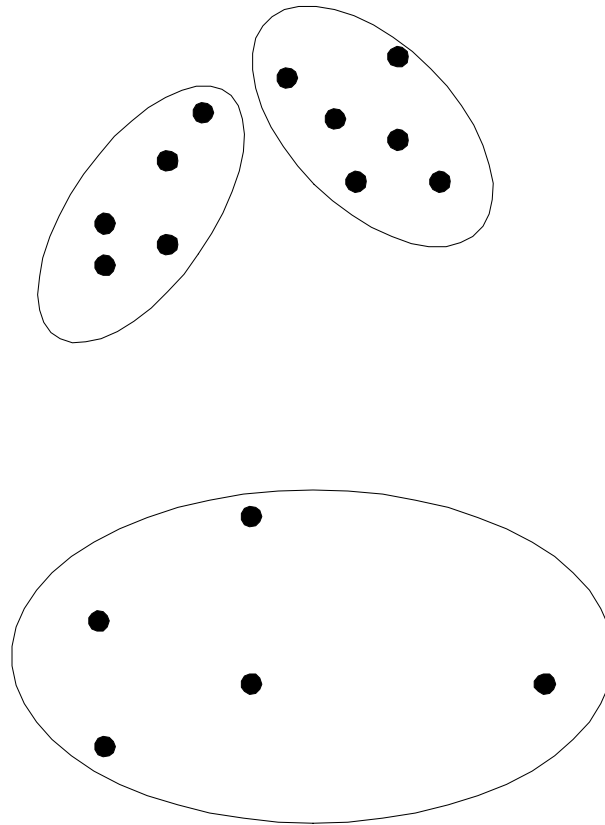
TYPES OF CLUSTERINGS

- 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

PARTITIONAL CLUSTERING

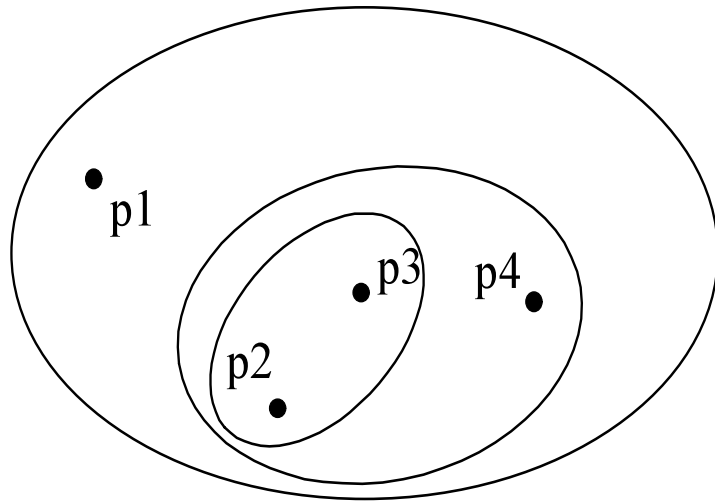


Original Points

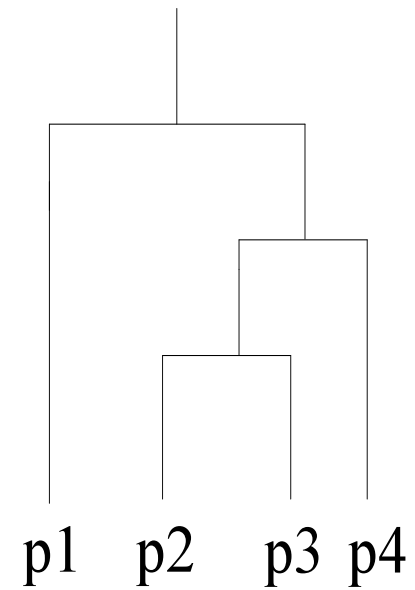


A Partitional Clustering

HIERARCHICAL CLUSTERING



Hierarchical Clustering



Dendrogram

CHARACTERISTICS OF THE INPUT 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

CLUSTERING ALGORITHMS

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

K-MEANS CLUSTERING

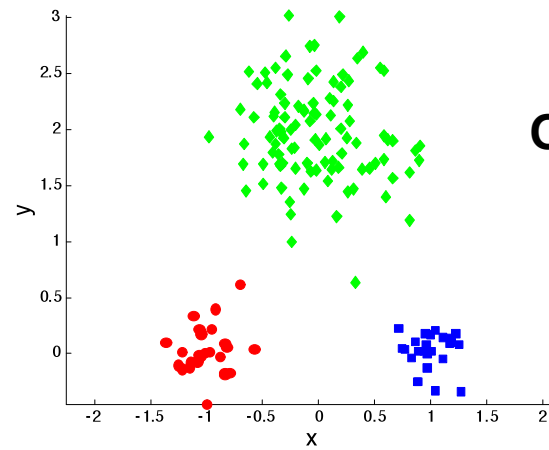
- Partitional clustering approach
- Each cluster is associated with a **centroid** (center point)
- Each point is assigned to the cluster with the closest centroid
- **Number of clusters, K , must be specified**
- The basic algorithm is very simple

-
- 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
-

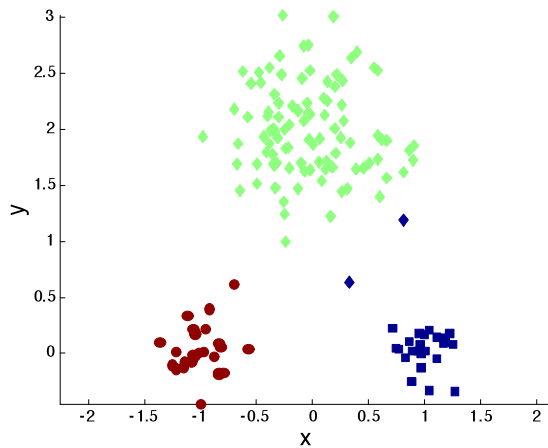
K-MEANS CLUSTERING – DETAILS

- 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’
- Complexity is $O(n * K * I * d)$
 - n = number of points, K = number of clusters,
 I = number of iterations, d = number of attributes

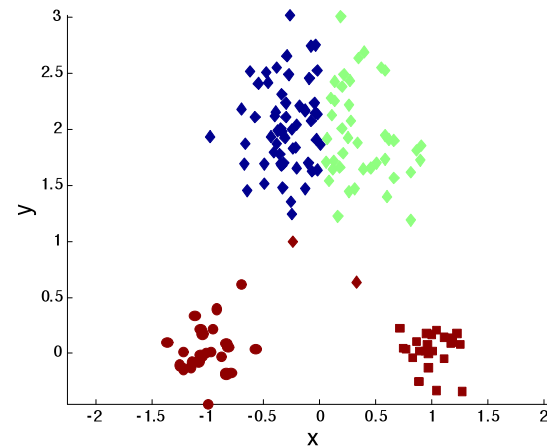
TWO DIFFERENT K-MEANS CLUSTERINGS



Original Points

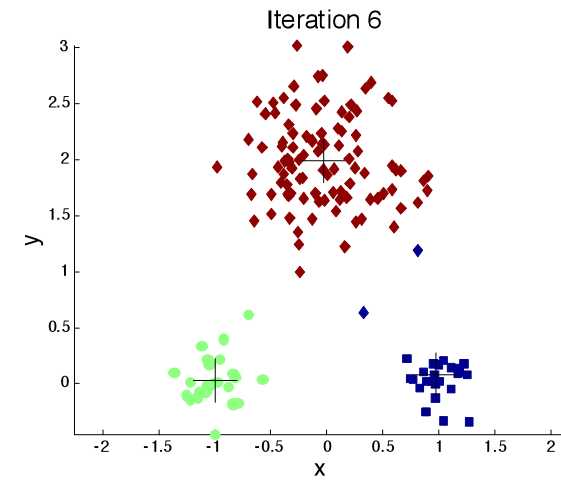
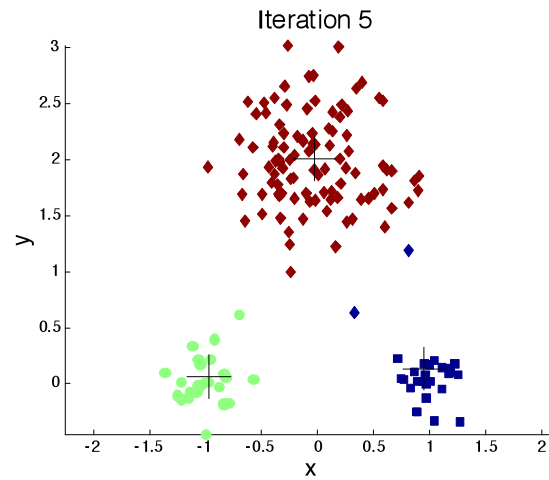
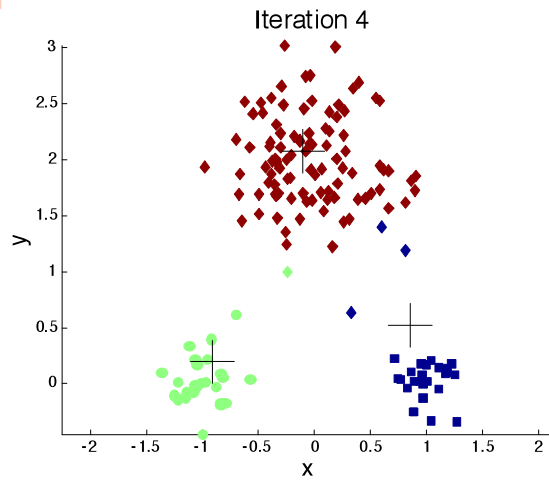
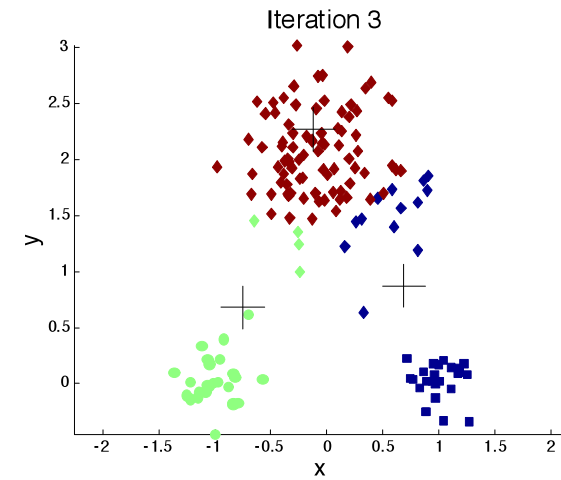
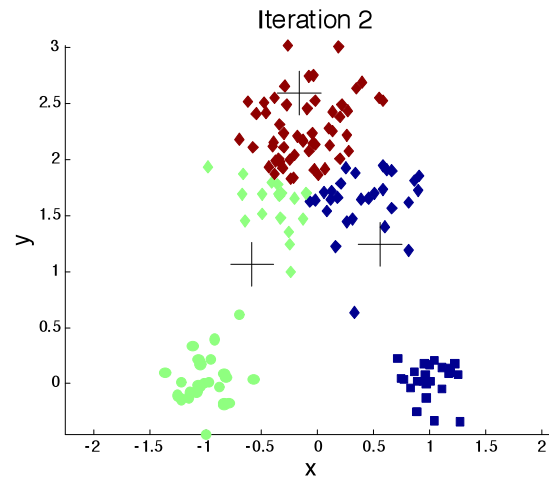
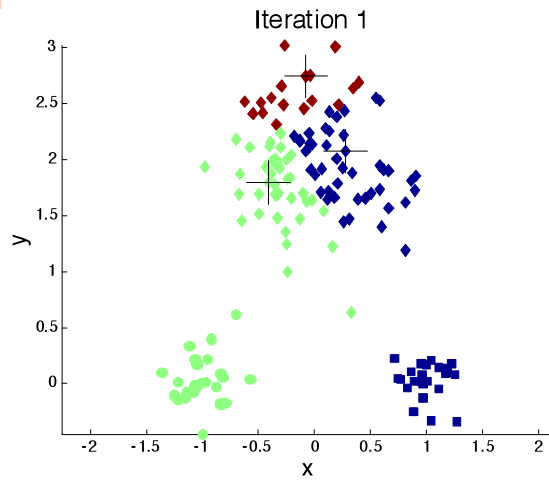


Optimal Clustering

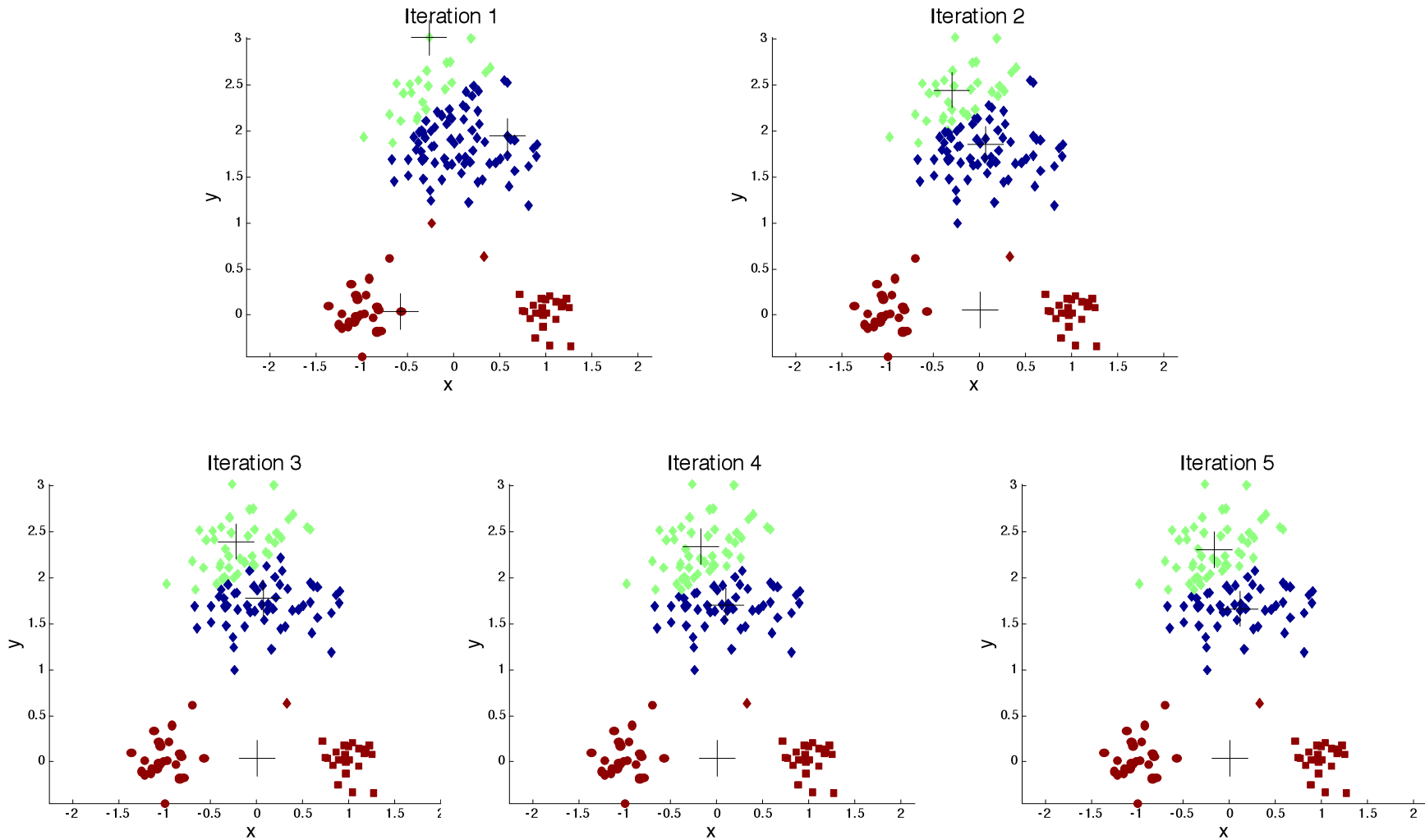


Sub-optimal Clustering

IMPORTANCE OF CHOOSING INITIAL CENTROIDS



IMPORTANCE OF CHOOSING INITIAL CENTROIDS ...

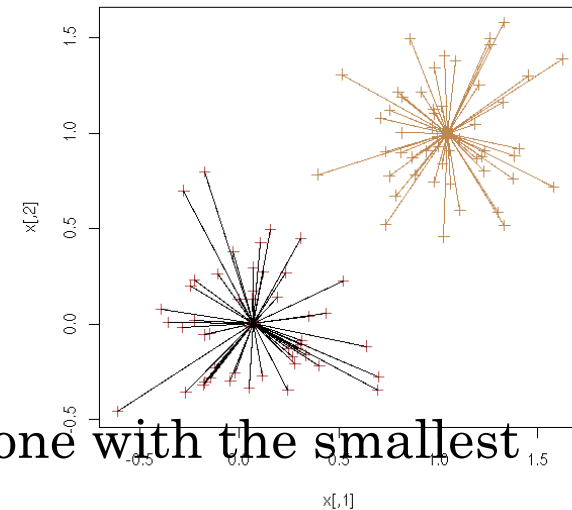


EVALUATING K-MEANS CLUSTERS

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

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

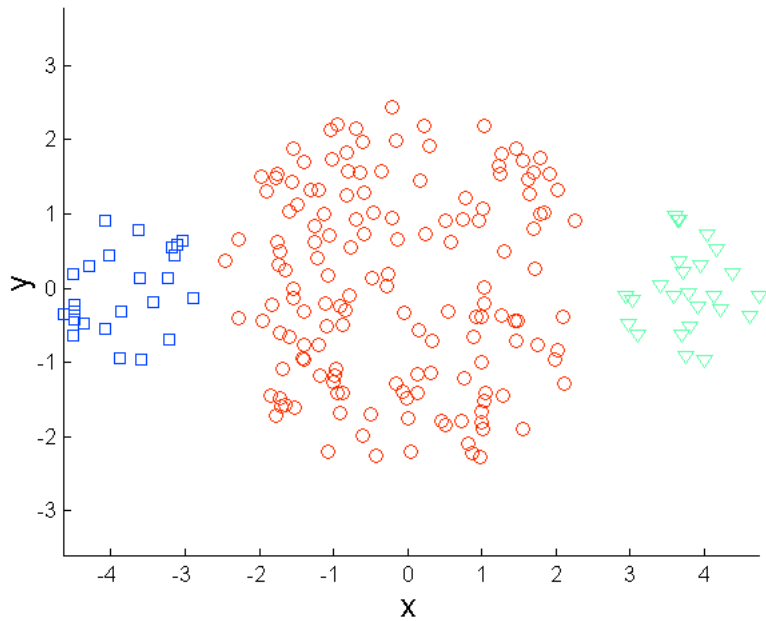
- x is a data point in cluster C_i and m_i is the representative point for cluster C_i
- Given two clusters, we can choose the one with the smallest error
- One easy way to reduce SSE is to increase K , the number of clusters
 - ◆ **A good clustering with smaller K can have a lower SSE than a poor clustering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



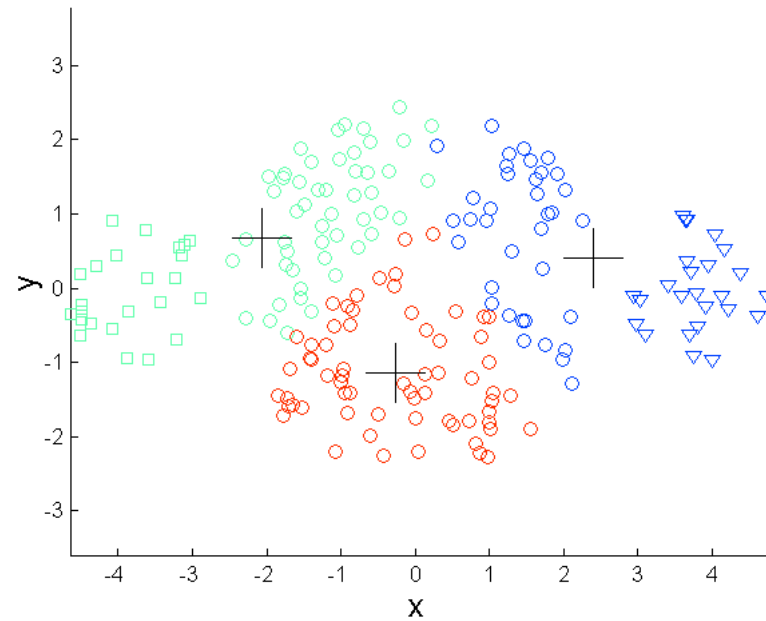
LIMITATIONS OF K-MEANS

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

LIMITATIONS OF K-MEANS: DIFFERING SIZES

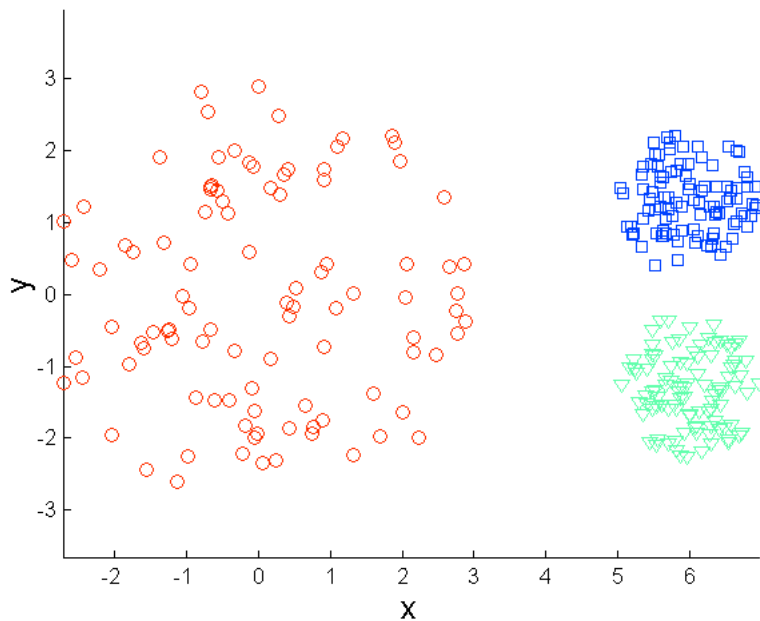


Original Points

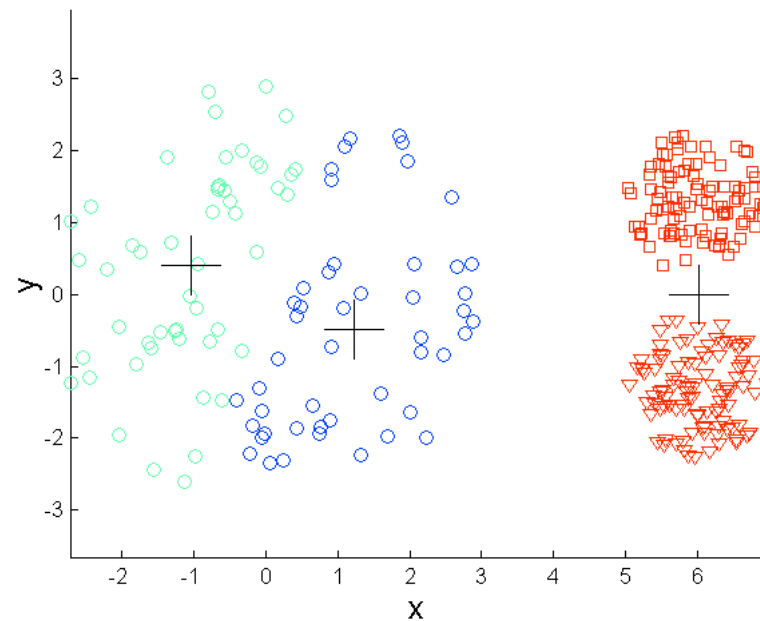


K-means (3 Clusters)

LIMITATIONS OF K-MEANS: DIFFERING DENSITY

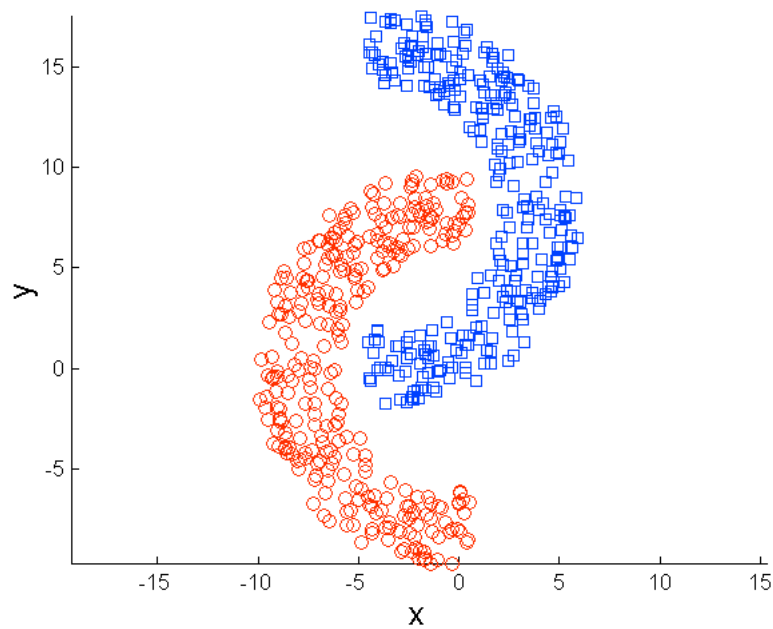


Original Points

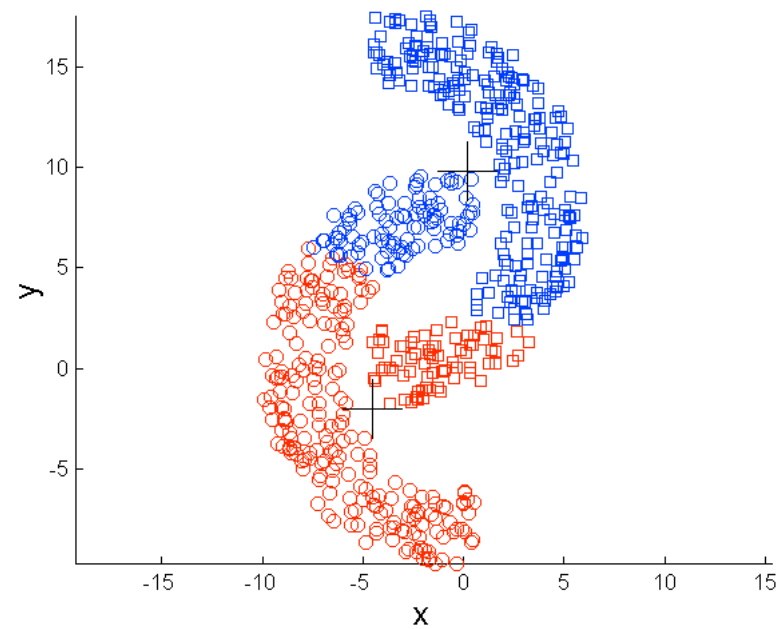


K-means (3 Clusters)

LIMITATIONS OF K-MEANS: NON-GLOBULAR SHAPES

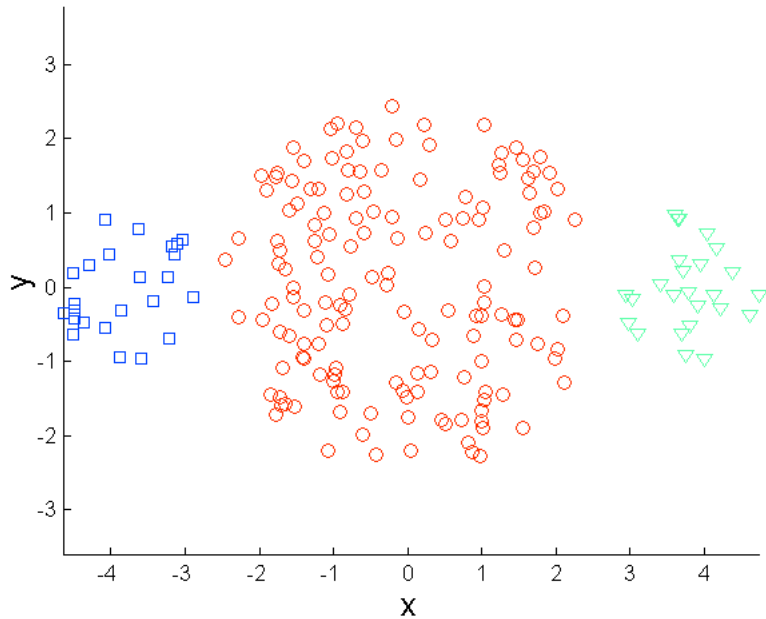


Original Points

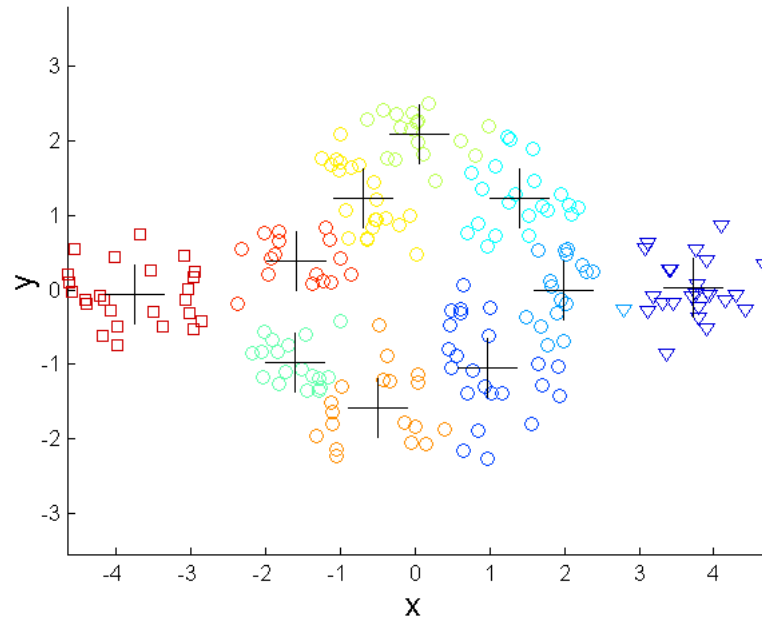


K-means (2 Clusters)

OVERCOMING K-MEANS LIMITATIONS



Original Points

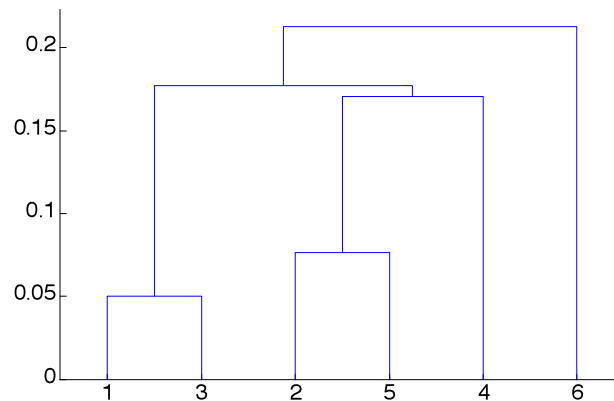


K-means Clusters

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

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



STRENGTHS OF HIERARCHICAL CLUSTERING

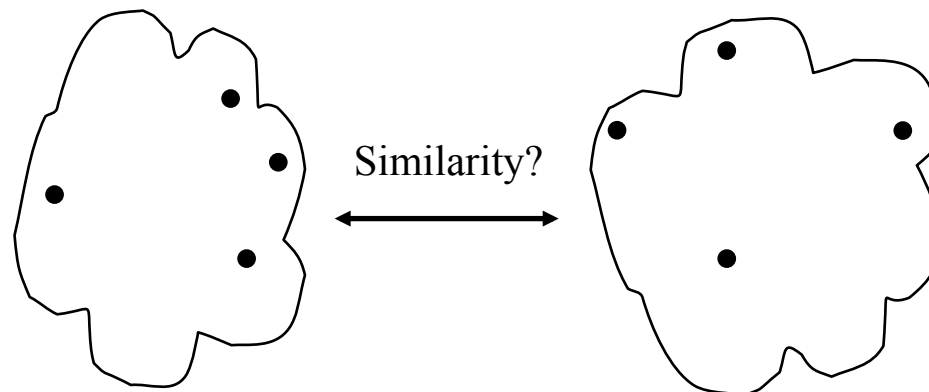
- Do not have to assume any particular number of clusters
 - Any desired number of clusters can be obtained by ‘cutting’ the dendrogram at the proper level
- They may correspond to meaningful taxonomies
 - Example in biological sciences (e.g., animal kingdom, phylogeny reconstruction, ...)

HIERARCHICAL CLUSTERING

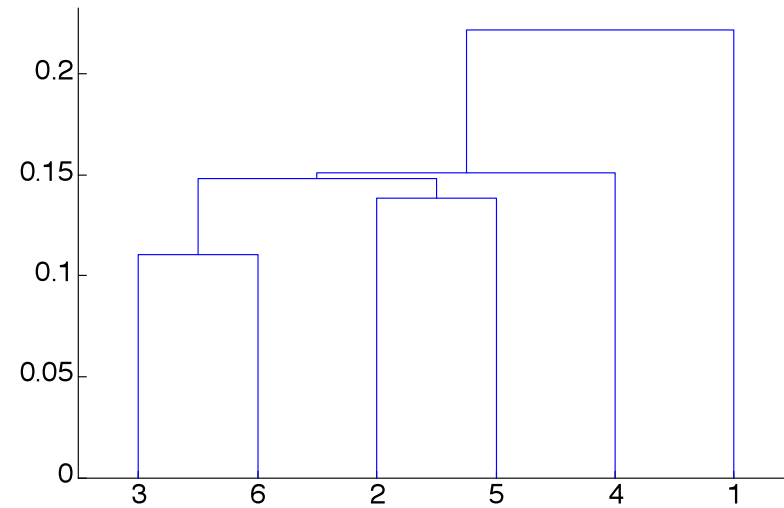
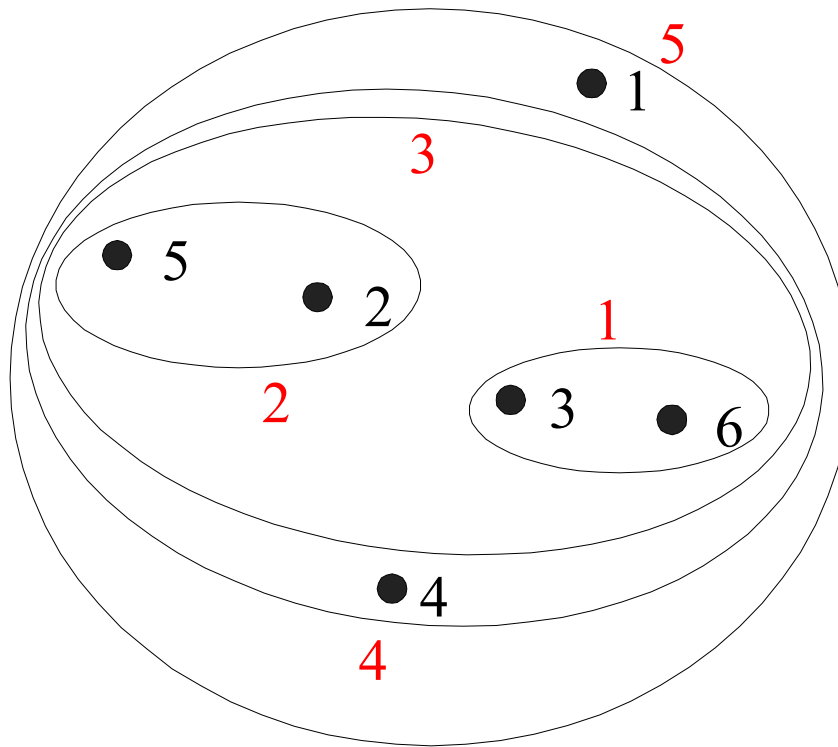
- Two main types of hierarchical clustering
 - Agglomerative:
 - ◆ Start with the points as individual clusters
 - ◆ At each step, merge the closest pair of clusters until only one cluster (or k clusters) left
 - Divisive:
 - ◆ Start with one, all-inclusive cluster
 - ◆ At each step, split a cluster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

AGGLOMERATIVE CLUSTERING ALGORITHM

- More popular hierarchical clustering technique
- Basic algorithm is straightforward
 - *Compute the proximity matrix*
 - *Let each data point be a cluster*
 - **Repeat**
 - *Merge the two closest clusters*
 - *Update the proximity matrix*
 - *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



HIERARCHICAL CLUSTERING: MIN



HIERARCHICAL CLUSTERING: TIME AND SPACE REQUIREMENTS

$O(N^2)$ space

$O(N^3)$ time in many cases

- There are N steps and at each step the proximity matrix (size: $O(N^2)$) must be updated and searched
- Complexity can be reduced to $O(N^2 \log(N))$ time for some approaches

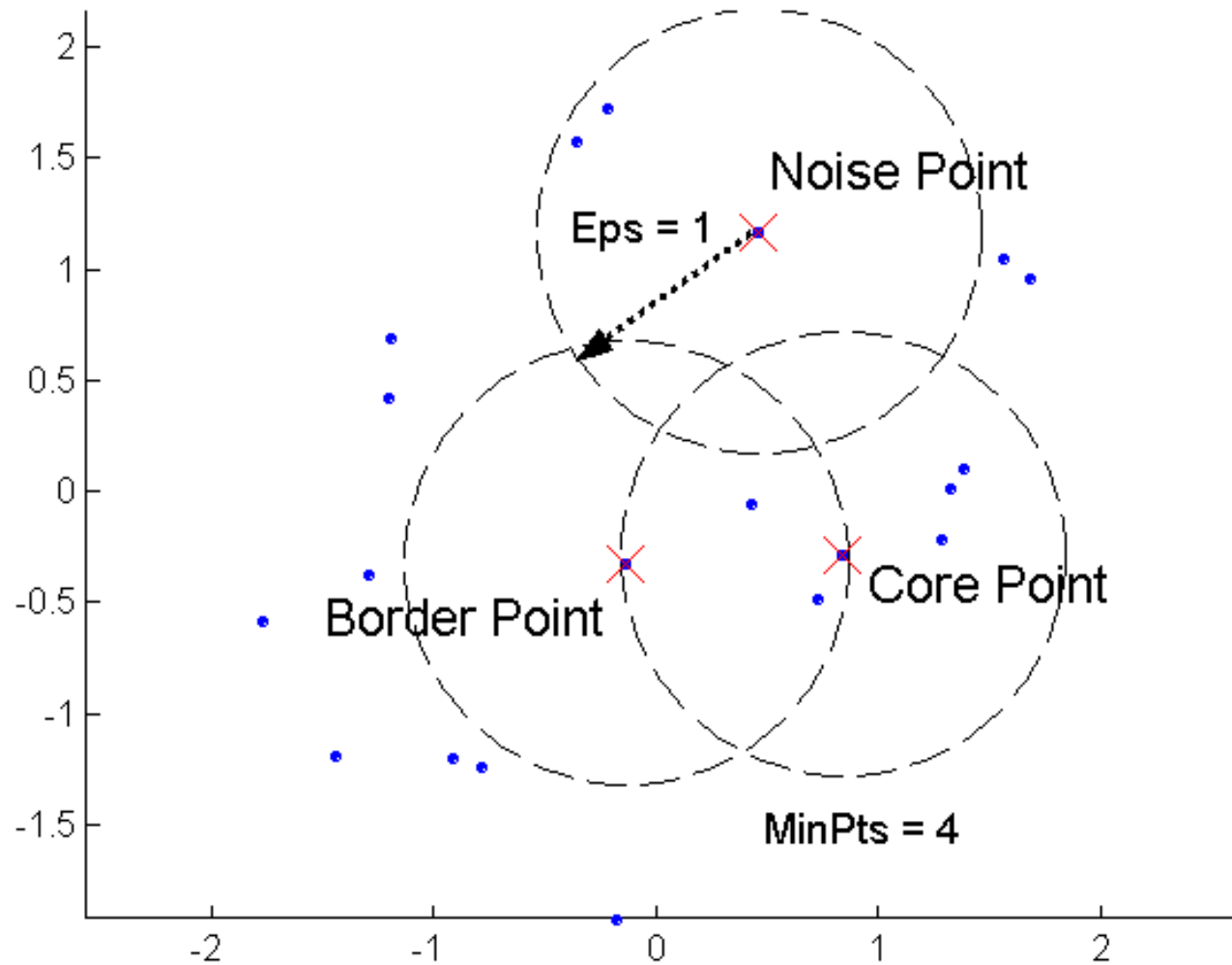
HIERARCHICAL CLUSTERING: PROBLEMS AND LIMITATIONS

- 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

DBSCAN

- DBSCAN is a density-based algorithm.
 - Density = number of points within a specified radius (Eps)
 - A point is a **core point** if it has more than a specified number of points (MinPts) within Eps
 - ◆ These are points that are at the interior of a cluster
 - A **border point** has fewer than MinPts within Eps, but is in the neighborhood of a core point
 - A **noise point** is any point that is not a core point or a border point.

DBSCAN: CORE, BORDER, AND NOISE POINTS



DBSCAN ALGORITHM

- Eliminate noise points
- Perform clustering on the remaining points

Algorithm 8.4 DBSCAN algorithm.

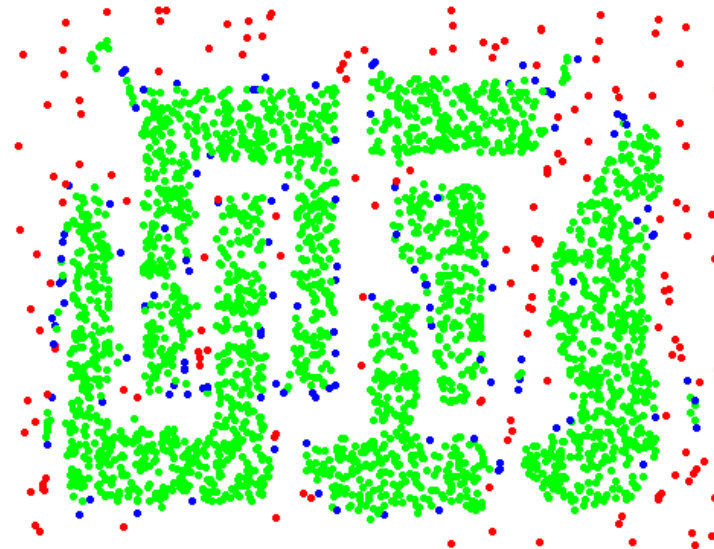
- 1: Label all points as core, border, or noise points.
 - 2: Eliminate noise points.
 - 3: Put an edge between all core points that are within Eps of each other.
 - 4: Make each group of connected core points into a separate cluster.
 - 5: Assign each border point to one of the clusters of its associated core points.
-

Complexity is $O(n^2)$ in the worst case. With low dimensionality and good data structure can reduce to $O(m \log m)$

DBSCAN: CORE, BORDER AND NOISE POINTS



Original Points



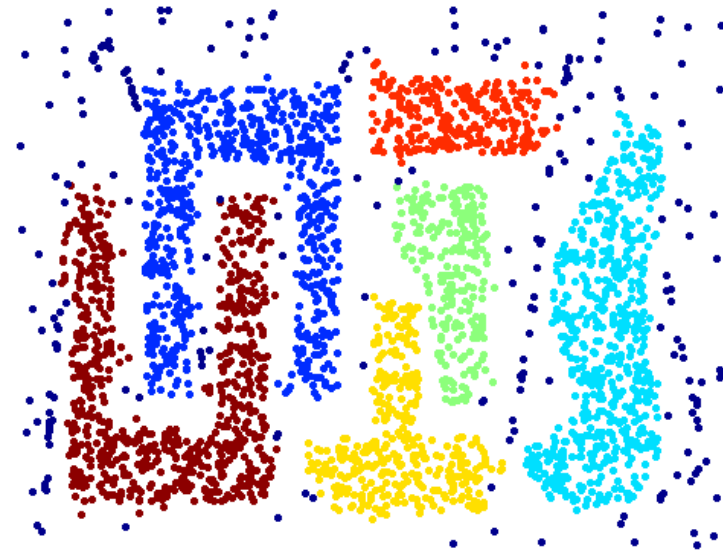
Point types: **core**,
border and **noise**

Eps = 10, MinPts = 4

WHEN DBSCAN WORKS WELL



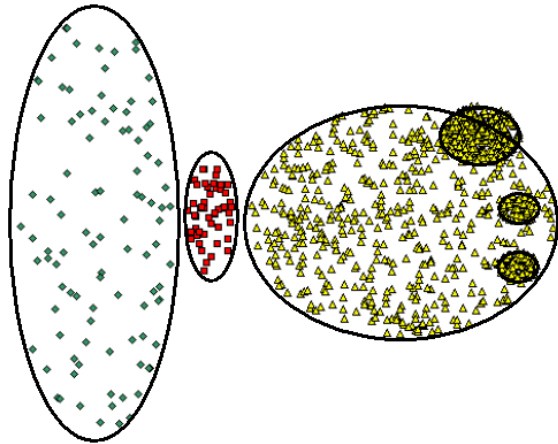
Original Points



Clusters

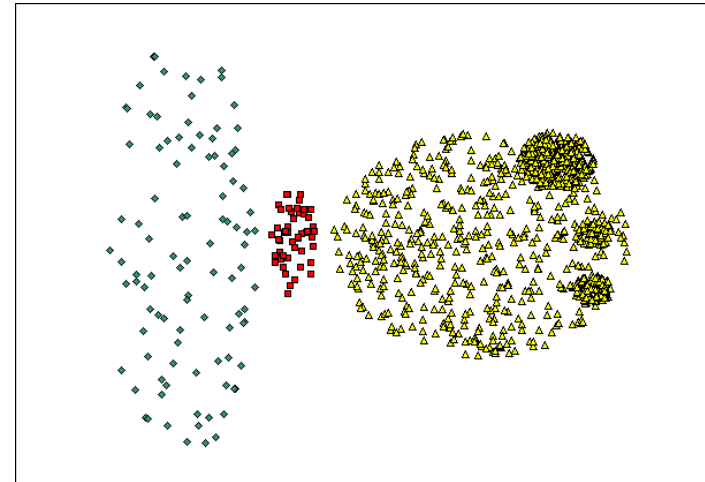
- **Resistant to Noise**
- **Can handle clusters of different shapes and sizes**

WHEN DBSCAN DOES NOT WORK WELL

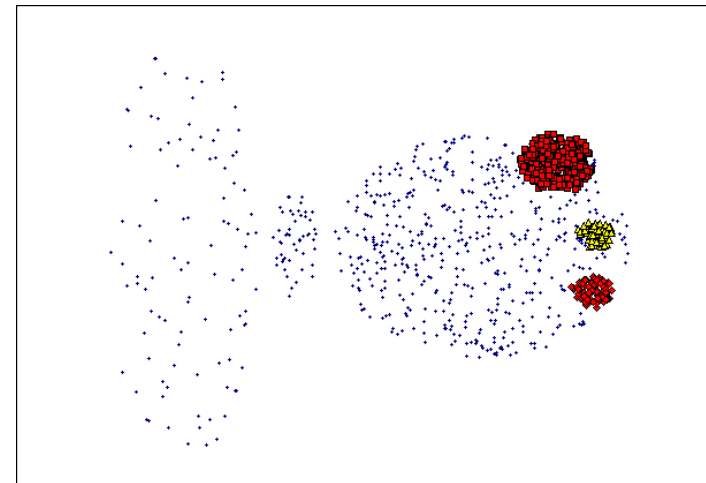


Original Points

- **Varying densities**
- **High-dimensional data**



(MinPts=4, Eps=9.92)



(MinPts=4, Eps=9.75).

CLUSTER VALIDITY

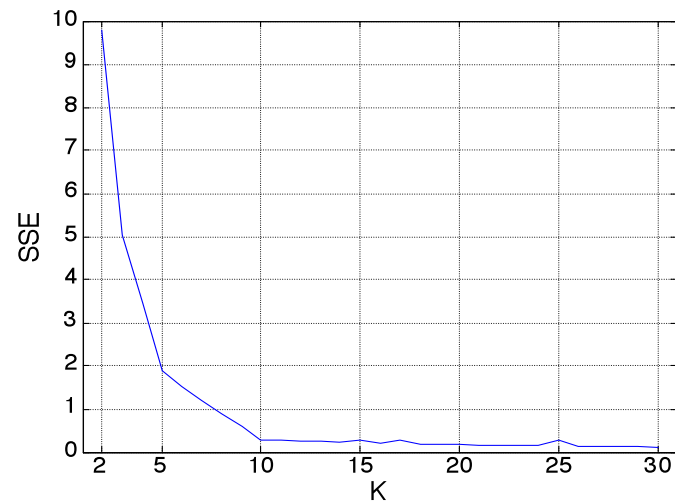
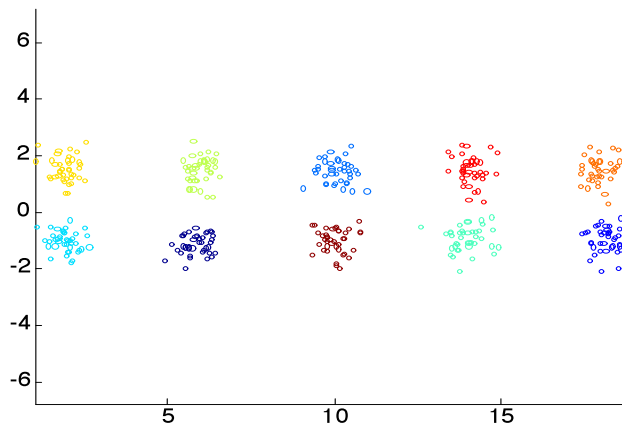
- For supervised classification we have a variety of measures to evaluate how good our model is
 - Accuracy, precision, recall
- For cluster analysis, the analogous question is how to evaluate the “goodness” of the resulting clusters?
- But “clusters are in the eye of the beholder”!
- Then why do we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

DIFFERENT ASPECTS OF CLUSTER VALIDATION

- Determining the **clustering tendency** of a set of data, i.e., distinguishing whether non-random structure actually exists in the data.
 1. Comparing the results of a cluster analysis to externally known results, e.g., to externally given class labels.
- Evaluating how well the results of a cluster analysis fit the data *without* reference to external information - only the data
 1. Comparing the results of two different sets of cluster analyses to determine which is better.
 2. Determining the 'correct' number of clusters.

INTERNAL MEASURES: SSE

- Clusters in more complicated figures aren't well separated
- Internal Index: Used to measure the goodness of a clustering structure without respect to external information
 - Sum of Square Error
- SSE is good for comparing two clusterings or two clusters (average SSE).
- Can also be used to estimate the number of clusters



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

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