

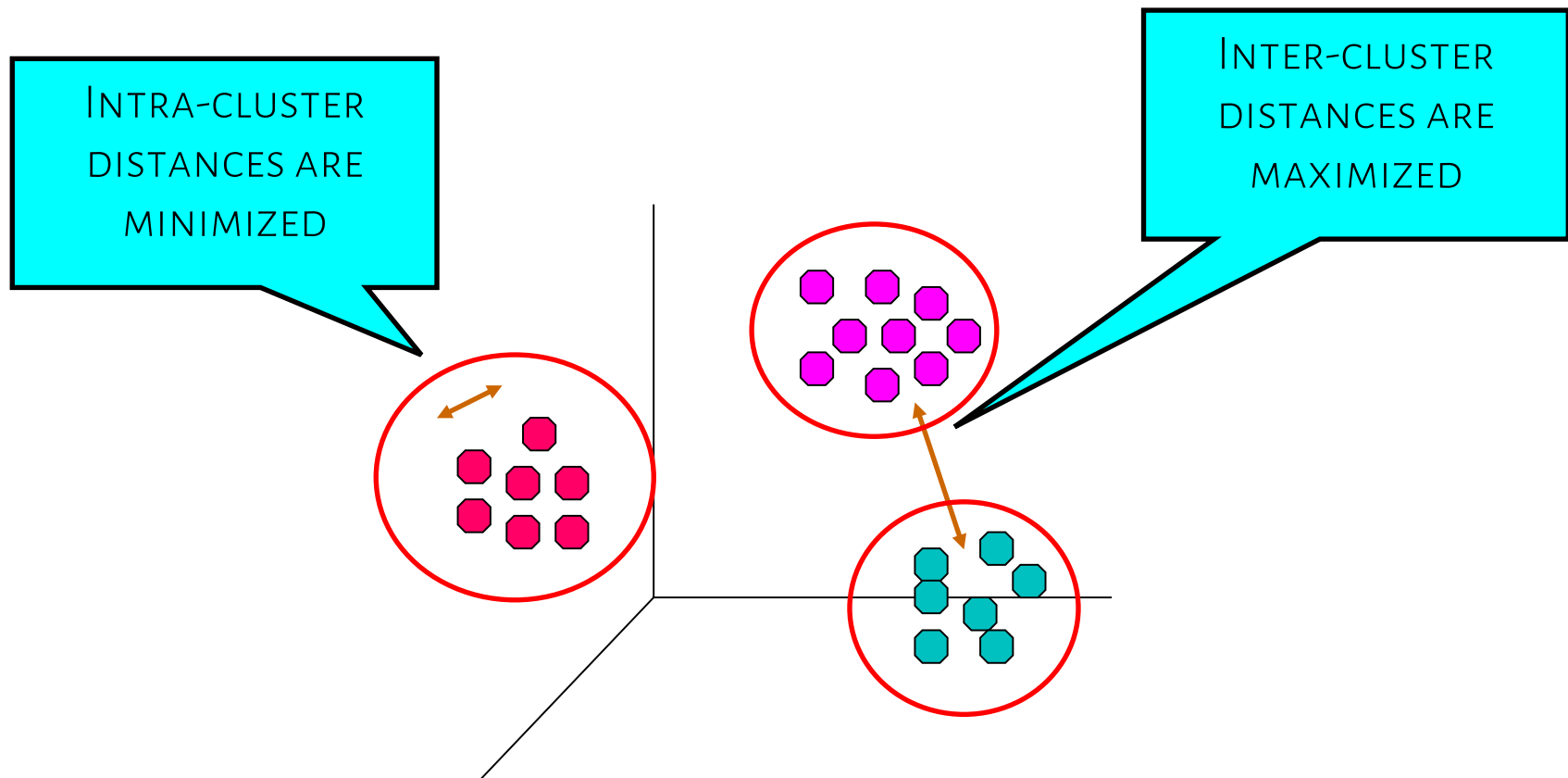
# Unit 5: Clustering

# Unit 5

## **Section 1: Concepts and algorithms**

# 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
- ✗ THE CONCEPT IS LESS “CLEAR” THAN CLASSIFICATION



# Intra-Cluster distances

- ✕ 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.

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

- ✕  $x$  is a data point in cluster  $C_i$  and  $m_i$  is the representative point for cluster  $C_i$
- ✕  $m_i$  usually corresponds to the center (mean) of the cluster

# Applications of Cluster Analysis

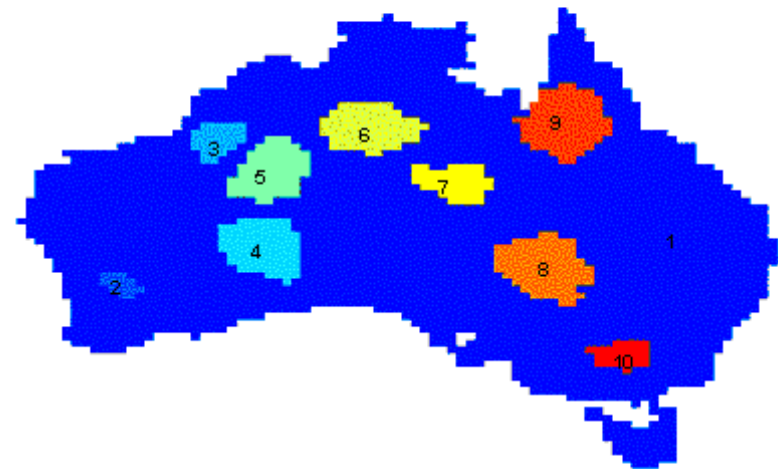
## ✕ UNDERSTANDING

✕ 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>
<b>1</b>	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, Oracl-DOWN, SGI-DOWN, Sun-DOWN	Technology1-DOWN
<b>2</b>	Apple-Comp-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	Technology2-DOWN
<b>3</b>	Fannie-Mae-DOWN, Fed-Home-Loan-DOWN, MBNA-Corp-DOWN, Morgan-Stanley-DOWN	Financial-DOWN
<b>4</b>	Baker-Hughes-UP, Dresser-Inds-UP, Halliburton-HLD-UP, Louisiana-Land-UP, Phillips-Petro-UP, Unocal-UP, Schlumberger-UP	Oil-UP

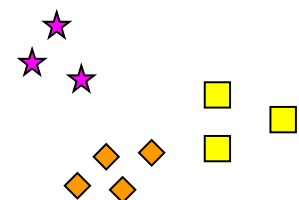
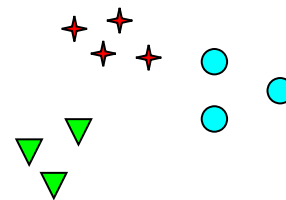
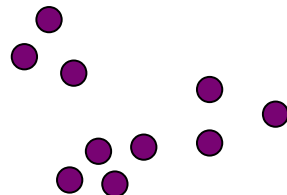
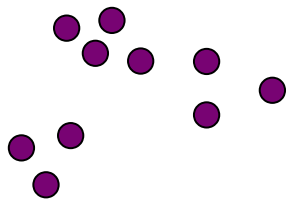
## ✕ SUMMARIZATION

✕ Reduce the size of large data sets



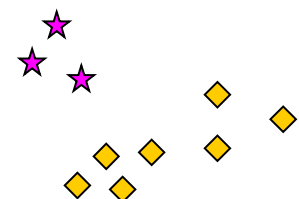
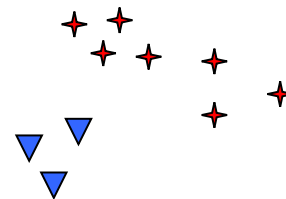
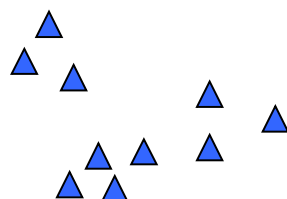
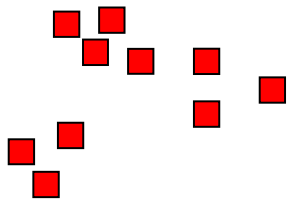
CLUSTERING PRECIPITATION IN AUSTRALIA

# Notion of a Cluster can be Ambiguous



HOW MANY CLUSTERS?

SIX CLUSTERS



TWO CLUSTERS

FOUR CLUSTERS

# Requirements of Clustering in Data Mining

- ✗ SCALABILITY
- ✗ ABILITY TO DEAL WITH DIFFERENT TYPES OF ATTRIBUTES
- ✗ ABILITY TO HANDLE DYNAMIC DATA
- ✗ DISCOVERY OF CLUSTERS WITH ARBITRARY SHAPE
- ✗ MINIMAL REQUIREMENTS FOR DOMAIN KNOWLEDGE TO DETERMINE INPUT PARAMETERS
- ✗ ABLE TO DEAL WITH NOISE AND OUTLIERS
- ✗ INSENSITIVE TO ORDER OF INPUT RECORDS
- ✗ HIGH DIMENSIONALITY
- ✗ INCORPORATION OF USER-SPECIFIED CONSTRAINTS
- ✗ INTERPRETABILITY AND USABILITY

# Most common Input Data Structures

- DATA MATRIX

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

- DISSIMILARITY MATRIX

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



# Type of data in clustering analysis

- ✕ INTERVAL-SCALED VARIABLES
- ✕ BINARY VARIABLES
- ✕ NOMINAL, ORDINAL, AND RATIO VARIABLES
- ✕ VARIABLES OF MIXED TYPES

# Interval-valued variables

## ✕STANDARDIZE DATA

✕Calculate the mean absolute deviation:

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

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

✕Calculate the standardized measurement (z-score)

✕USING MEAN ABSOLUTE DEVIATION IS MORE ROBUST THAN USING STANDARD DEVIATION

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

# Similarity and Dissimilarity Between Objects

✕ DISTANCES ARE NORMALLY USED TO MEASURE THE SIMILARITY OR DISSIMILARITY BETWEEN TWO DATA OBJECTS

✕ SOME POPULAR ONES INCLUDE: MINKOWSKI DISTANCE:

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

✕ where  $i = (x_{i1}, x_{i2}, \dots, x_{ip})$  and  $j = (x_{j1}, x_{j2}, \dots, x_{jp})$  are two p-dimensional data objects, and q is a positive integer

✕ IF  $q = 1$ , D IS MANHATTAN DISTANCE

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

# Similarity and Dissimilarity Between Objects (Cont.)

✕ IF  $Q = 2$ ,  $D$  IS EUCLIDEAN DISTANCE:

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

✕ Properties

1)  $d(i, j) \geq 0$

2)  $d(i, i) = 0$

3)  $d(i, j) = d(j, i)$

4)  $d(i, j) \leq d(i, k) + d(k, j)$

✕ ALSO, ONE CAN USE

✕ weighted distance

✕ parametric Pearson product moment correlation

✕ other dissimilarity measures

# Binary Variables

## ✕ A CONTINGENCY TABLE FOR BINARY DATA

		Object $j$		
		1	0	$sum$
Object $i$	1	$a$	$b$	$a+b$
	0	$c$	$d$	$c+d$
$sum$		$a+c$	$b+d$	$p$

## ✕ DISTANCE MEASURE FOR SYMMETRIC BINARY VARIABLES:

$$d(i,j) = \frac{b+c}{a+b+c+d}$$

## ✕ DISTANCE MEASURE FOR ASYMMETRIC BINARY VARIABLES:

$$d(i,j) = \frac{b+c}{a+b+c}$$

## ✕ JACCARD COEFFICIENT (SIMILARITY METRIC FOR ASYMMETRIC BINARY VARIABLES):

$$sim_{Jaccard}(i,j) = \frac{a}{a+b+c}$$

# Dissimilarity between Binary Variables

## ■ EXAMPLE

Name	Gender	Fever	Cough	Test-1	Test-2	Test-3	Test-4
Jack	M	Y	N	P	N	N	N
Mary	F	Y	N	P	N	P	N
Jim	M	Y	P	N	N	N	N

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

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

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

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

# Nominal Variables

- ✗ A GENERALIZATION OF THE BINARY VARIABLE IN THAT IT CAN TAKE MORE THAN 2 STATES, E.G., RED, YELLOW, BLUE, GREEN

- ✗ METHOD 1: SIMPLE MATCHING (EQUIVALENT TO HAMMING DISTANCE FOR BINARY VARIABLES)

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

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

- ✗ METHOD 2: USE A LARGE NUMBER OF BINARY VARIABLES

  - ✗ creating a new binary variable for each of the  $M$  nominal states

# Ordinal Variables

- ✗ AN ORDINAL VARIABLE CAN BE DISCRETE OR CONTINUOUS
- ✗ ORDER IS IMPORTANT, E.G., RANK
- ✗ CAN BE TREATED LIKE INTERVAL-SCALED
  - ✗ replace  $x_{if}$  by their rank  $r_{if} \in \{1, \dots, M_f\}$
  - ✗ map the range of each variable onto  $[0, 1]$  by replacing  $i$ -th object in the  $f$ -th variable by

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

- ✗ compute the dissimilarity using methods for interval-scaled variables



# Ratio-Scaled Variables

- RATIO-SCALED VARIABLE: A POSITIVE MEASUREMENT ON A NONLINEAR SCALE, APPROXIMATELY AT EXPONENTIAL SCALE, SUCH AS  $AE^{B_T}$  OR  $AE^{-B_T}$
- METHODS:
  - treat them like interval-scaled variables—*not a good choice!* (why?—the scale can be distorted)
  - apply logarithmic transformation

$$Y_{IF} = \text{LOG}(X_{IF})$$

- treat them as continuous ordinal data treat their rank as interval-scaled

# Variables of Mixed Types

- A DATABASE MAY CONTAIN ALL THE SIX TYPES OF VARIABLES
  - symmetric binary, asymmetric binary, nominal, ordinal, interval and ratio
- ONE MAY USE A WEIGHTED FORMULA TO COMBINE THEIR EFFECTS

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

- $f$  is binary or nominal:
  - $d_{ij}^{(f)} = 0$  if  $x_{if} = x_{jf}$ , or  $d_{ij}^{(f)} = 1$  otherwise
- $f$  is interval-based: use the normalized distance
- $f$  is ordinal or ratio-scaled
  - compute ranks  $r_{if}$  and
  - and treat  $z_{if}$  as interval-scaled  $z_{if} = \frac{r_{if} - 1}{M_f - 1}$

# Vector Objects

- VECTOR OBJECTS: KEYWORDS IN DOCUMENTS, GENE FEATURES IN MICRO-ARRAYS, ETC.
- BROAD APPLICATIONS: INFORMATION RETRIEVAL, BIOLOGIC TAXONOMY, ETC.
- COSINE MEASURE

$$s(\vec{X}, \vec{Y}) = \frac{\vec{X}^t \cdot \vec{Y}}{|\vec{X}| |\vec{Y}|},$$

- A VARIANT: TANIMOTO COEFFICIENT

$\vec{X}^t$  is a transposition of vector  $\vec{X}$ ,  $|\vec{X}|$  is the Euclidean normal of vector  $\vec{X}$ ,

$$s(\vec{X}, \vec{Y}) = \frac{\vec{X}^t \cdot \vec{Y}}{\vec{X}^t \cdot \vec{X} + \vec{Y}^t \cdot \vec{Y} - \vec{X}^t \cdot \vec{Y}},$$

# Types of Basic Clusterings

- ✕ A CLUSTERING IS A SET OF CLUSTERS
- ✕ TWO MAJOR TYPES:
  - ✕ Partitional clustering
  - ✕ Hierarchical clustering
- ✕ 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

# Other Clustering Approaches

## ✕PARTITIONING APPROACH:

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

## ✕HIERARCHICAL APPROACH:

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

## ✕DENSITY-BASED APPROACH:

- ✕Based on connectivity and density functions
- ✕Typical methods: DBSACN, OPTICS, DenClue

# Major Clustering Approaches (II)

## ✕GRID-BASED APPROACH:

- ✕based on a multiple-level granularity structure
- ✕Typical methods: STING, WaveCluster, CLIQUE

## ✕MODEL-BASED:

- ✕A model is hypothesized for each of the clusters and tries to find the best fit of that model to each other
- ✕Typical methods: EM, SOM, COBWEB

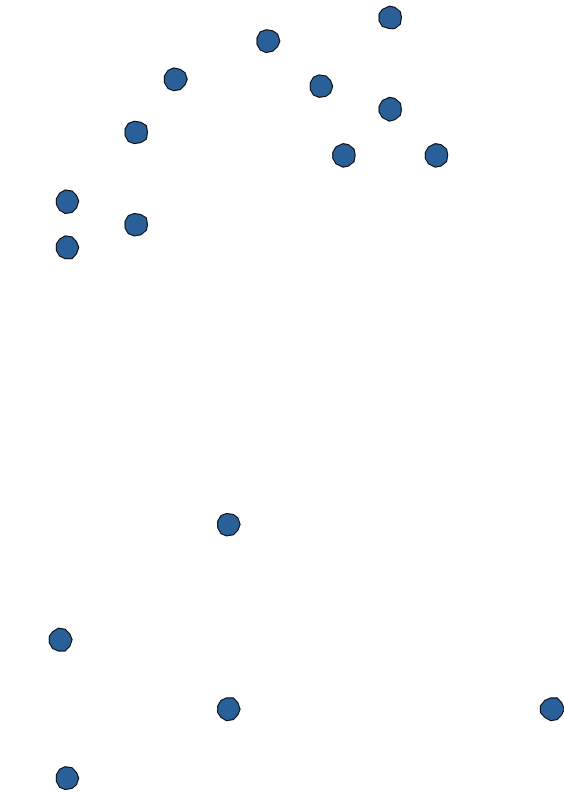
## ✕FREQUENT PATTERN-BASED:

- ✕Based on the analysis of frequent patterns
- ✕Typical methods: pCluster

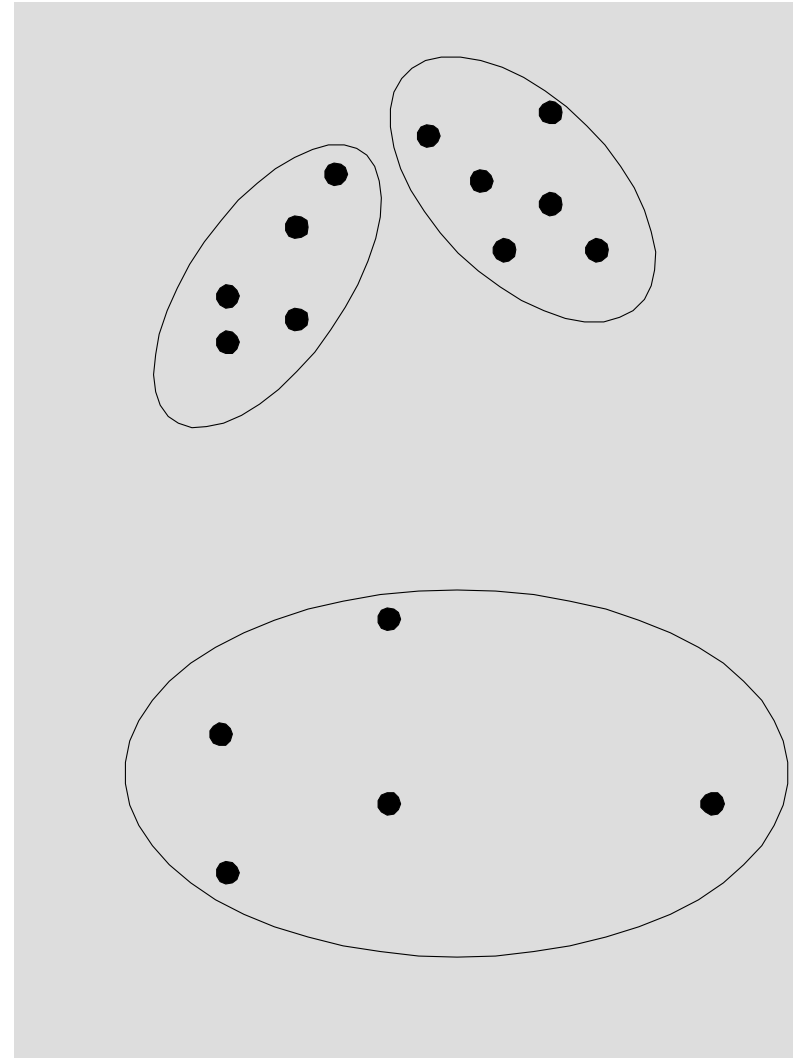
## ✕USER-GUIDED OR CONSTRAINT-BASED:

- ✕Clustering by considering user-specified or application-specific constraints
- ✕Typical methods: COD (obstacles), constrained clustering

# Partitional Clustering

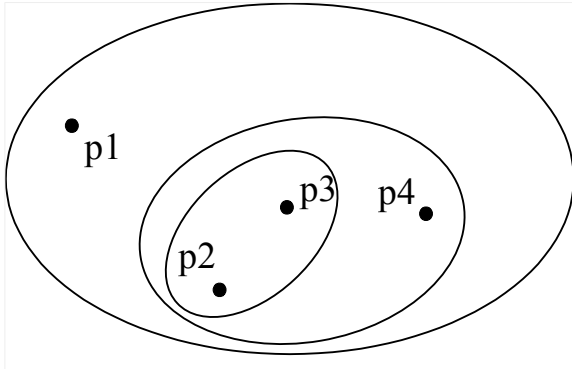


ORIGINAL POINTS

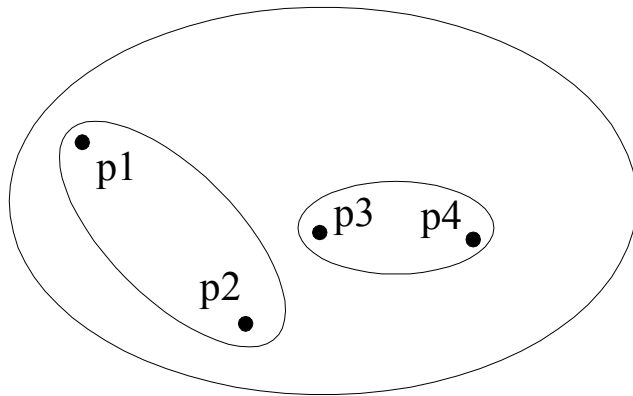


A PARTITIONAL CLUSTERING

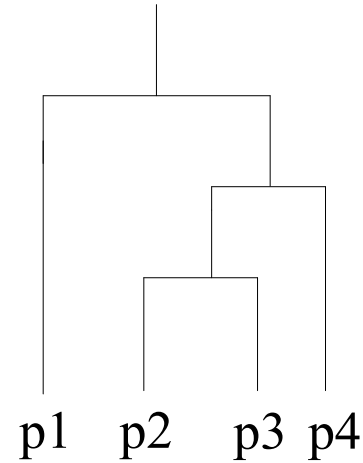
# Hierarchical Clustering



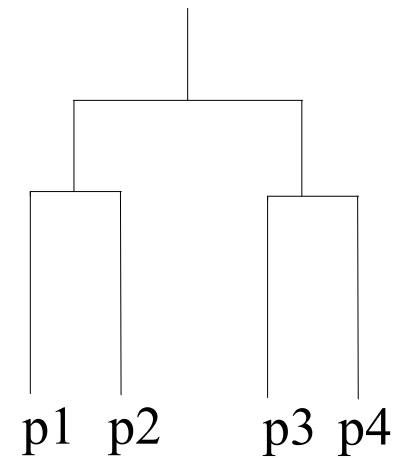
TRADITIONAL HIERARCHICAL CLUSTERING



NON-TRADITIONAL HIERARCHICAL CLUSTERING



TRADITIONAL DENDROGRAM



NON-TRADITIONAL DENDROGRAM



# Other Distinctions Between Sets of Clusters

## ✕EXCLUSIVE VERSUS NON-EXCLUSIVE

- ✕In non-exclusive clustering, points may belong to multiple clusters.
- ✕Can represent multiple classes or 'border' points

## ✕FUZZY VERSUS 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 VERSUS COMPLETE

- ✕In some cases, we only want to cluster some of the data

## ✕HETEROGENEOUS VERSUS HOMOGENEOUS

- ✕Cluster of widely different sizes, shapes, and densities

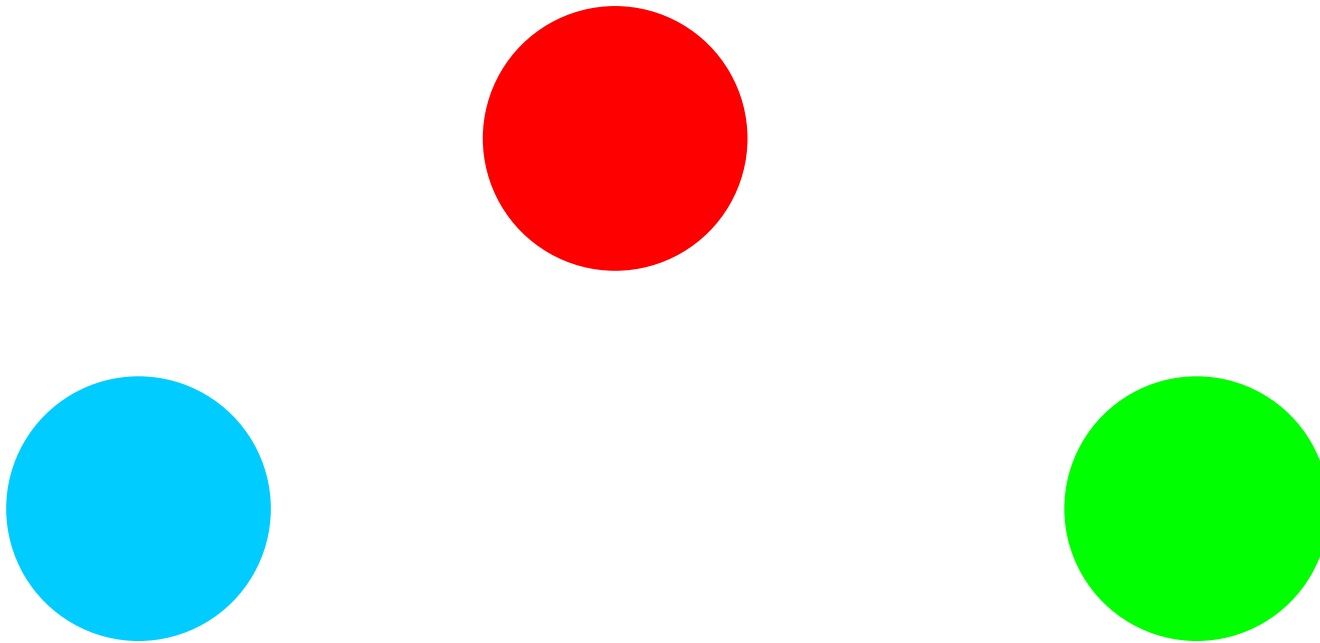
# Types of Clusters

- x WELL-SEPARATED CLUSTERS
- x CENTER-BASED CLUSTERS
- x CONTIGUOUS CLUSTERS
- x DENSITY-BASED CLUSTERS
- x PROPERTY OR CONCEPTUAL
- x DESCRIBED BY AN OBJECTIVE FUNCTION

# Types of Clusters: Well-Separated

## ✕WELL-SEPARATED CLUSTERS:

✕A cluster is a set of points such that any point in a cluster is closer (or more similar) to every other point in the cluster than to any point not in the cluster.

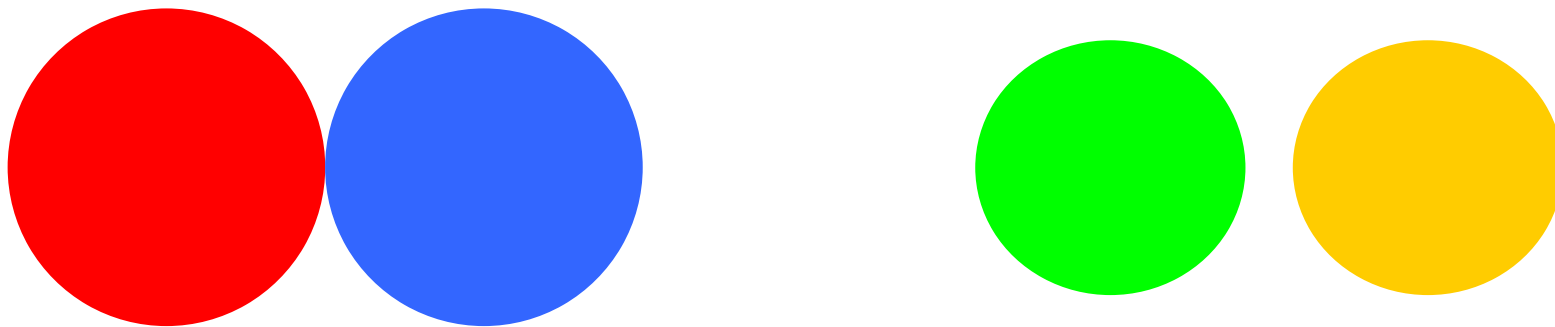


3 WELL-SEPARATED CLUSTERS

# Types of Clusters: Center-Based

## ✕CENTER-BASED

- ✕ A cluster is a set of objects such that an object in a cluster is closer (more similar) to the “center” of a cluster, than to the center of any other cluster
- ✕ The center of a cluster is often a centroid, the average of all the points in the cluster, or a medoid, the most “representative” point of a cluster

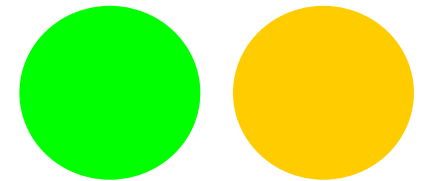
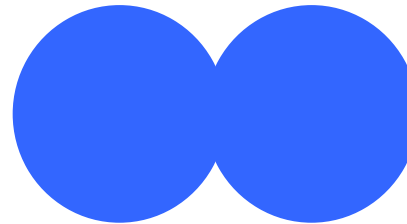
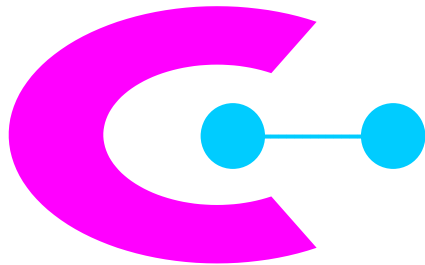
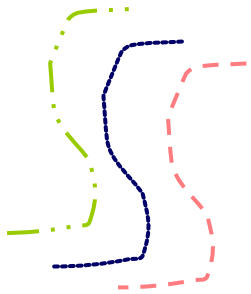


4 CENTER-BASED CLUSTERS

# Types of Clusters: Contiguity-Based

## ✕CONTIGUOUS CLUSTER (NEAREST NEIGHBOR OR TRANSITIVE)

✕A cluster is a set of points such that a point in a cluster is closer (or more similar) to one or more other points in the cluster than to any point not in the cluster.

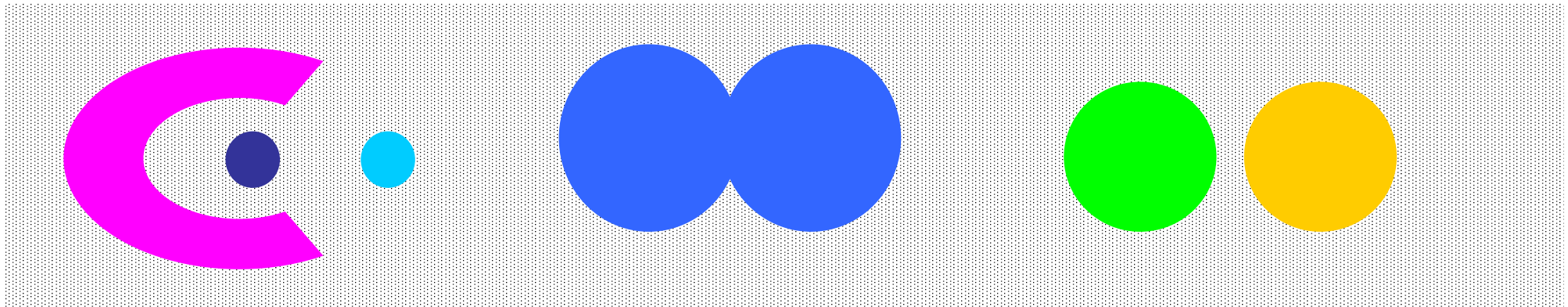


8 CONTIGUOUS CLUSTERS

# 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.

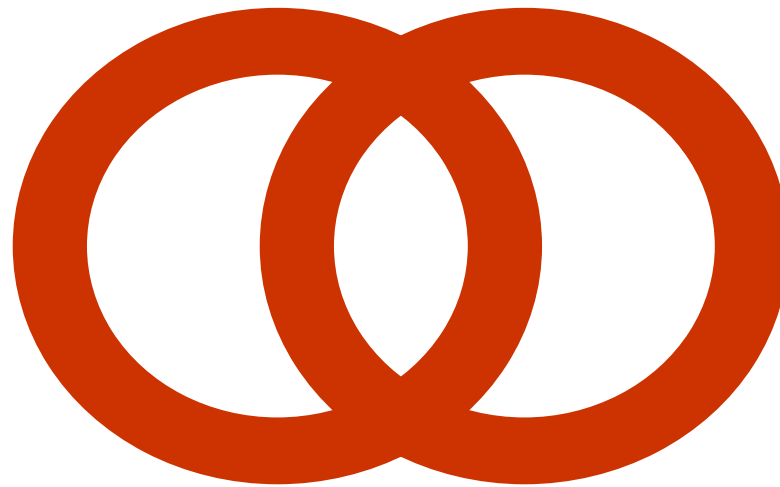


6 DENSITY-BASED CLUSTERS

# 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

# 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.



# Types of Clusters: Objective Function

- ✗ 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

# 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

# Partitional clustering

# K-means Clustering

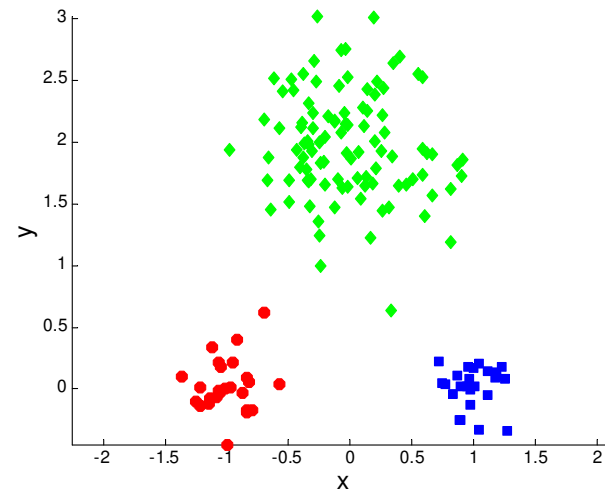
- ✕ PARTITIONAL CLUSTERING APPROACH
- ✕ EACH CLUSTER IS ASSOCIATED WITH A CENTROID (CENTER POINT)
- ✕ EACH POINT IS ASSIGNED TO THE CLUSTER WITH THE CLOSEST CENTROID
- ✕ NUMBER OF CLUSTERS,  $K$ , MUST BE SPECIFIED
- ✕ THE BASIC ALGORITHM IS VERY SIMPLE

- 
- 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
-

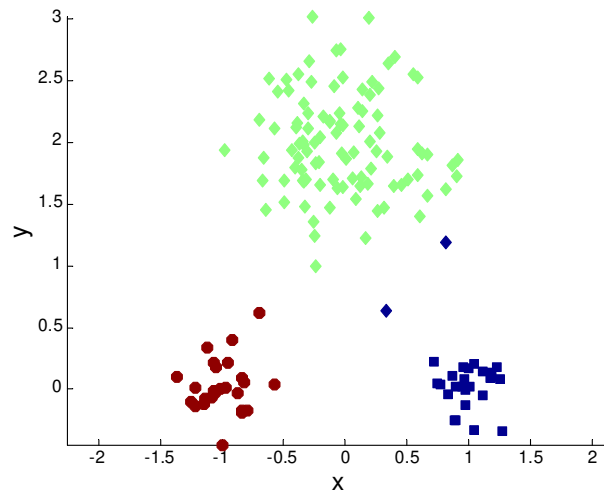
# K-means Clustering – Details

- ✗ 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  $O(N \times K \times I \times D)$ 
  - ✗  $n$  = number of points,  $K$  = number of clusters,  
 $I$  = number of iterations,  $d$  = number of attributes

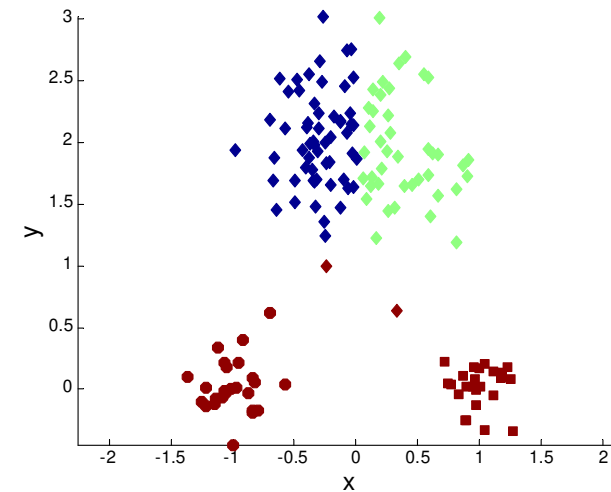
# Two different K-means Clusterings



ORIGINAL POINTS

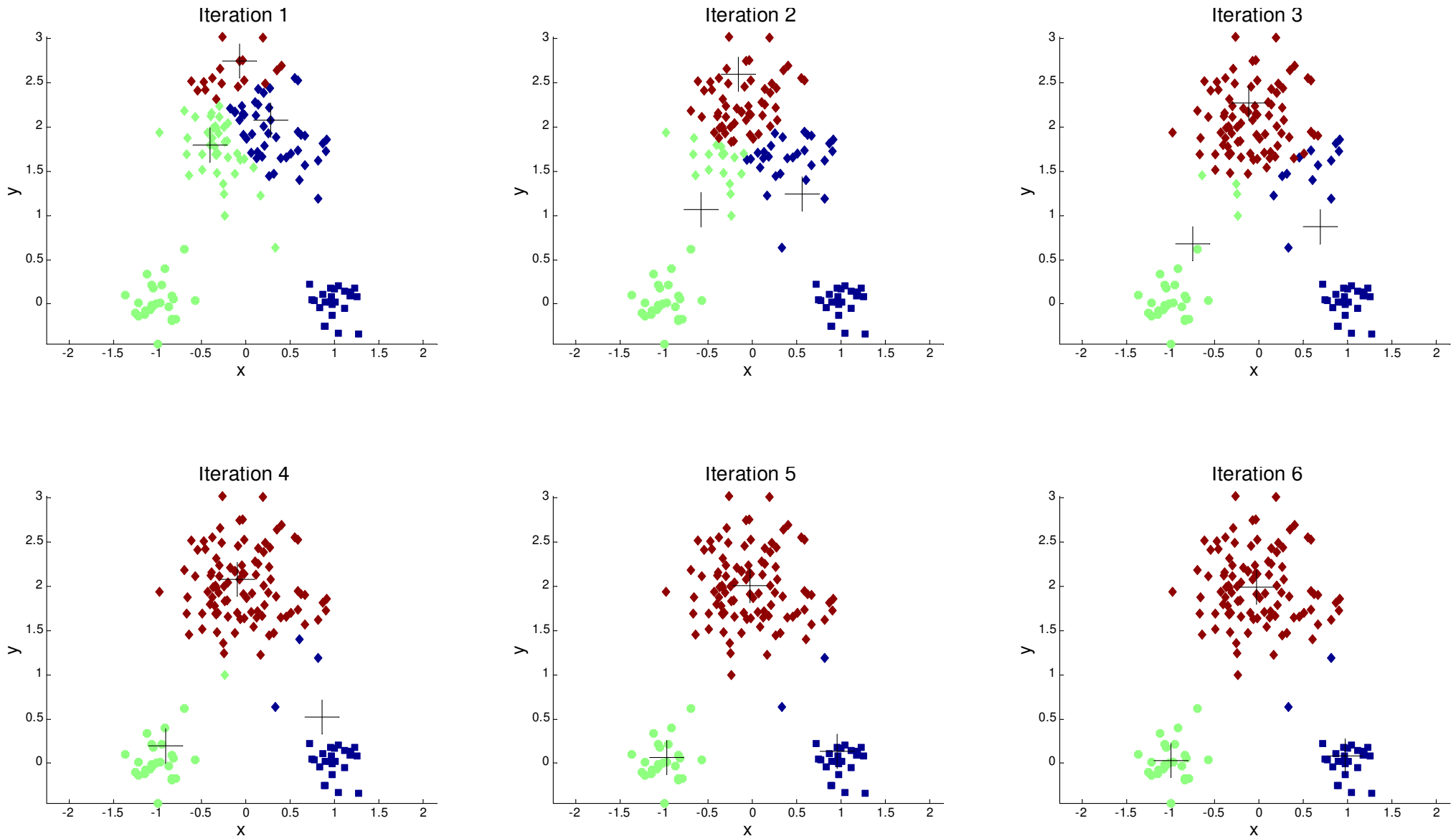


OPTIMAL CLUSTERING

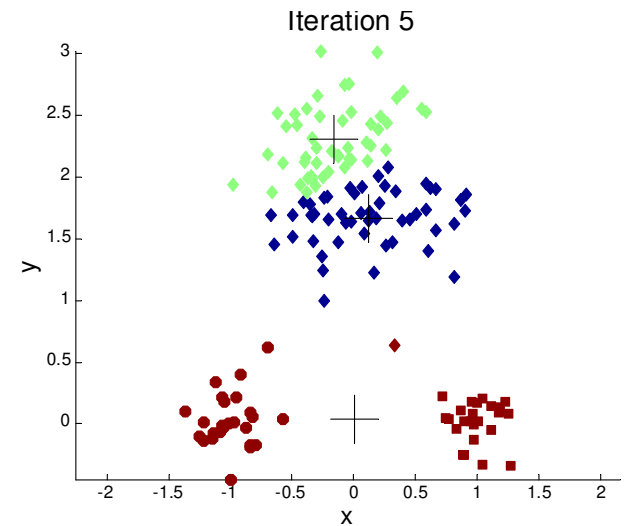
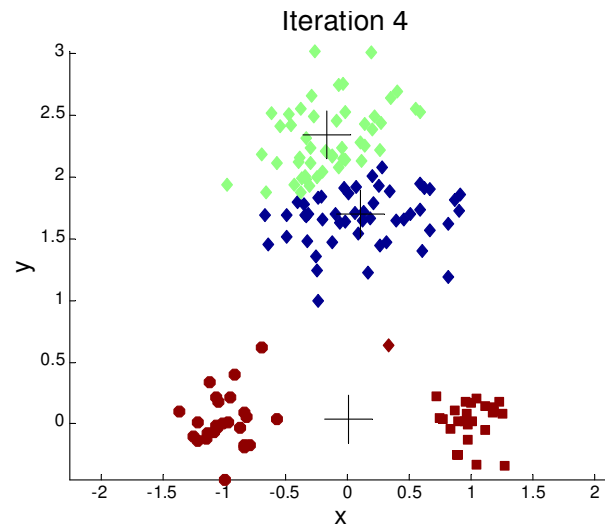
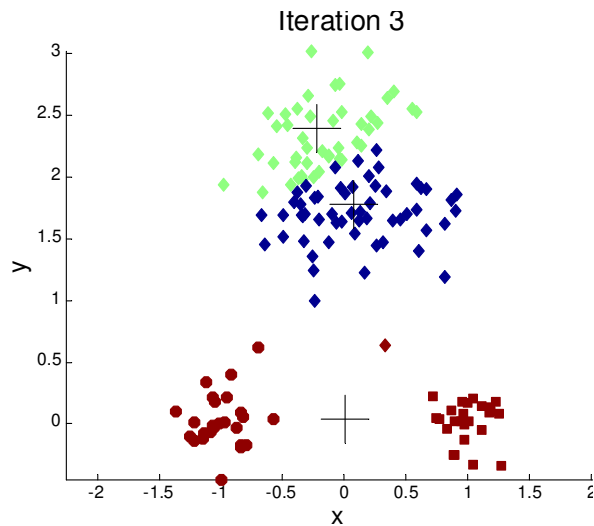
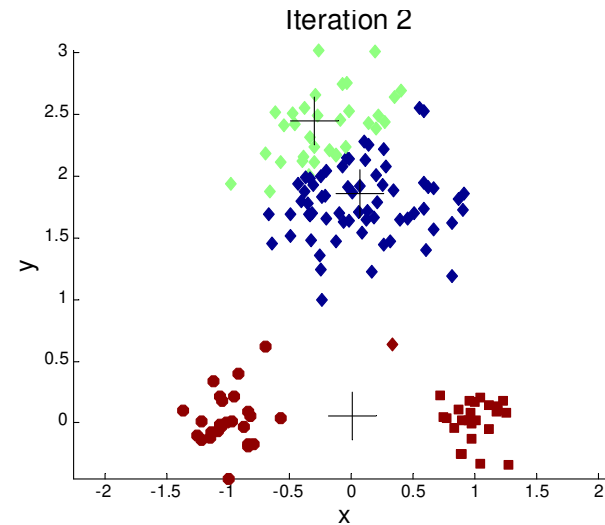
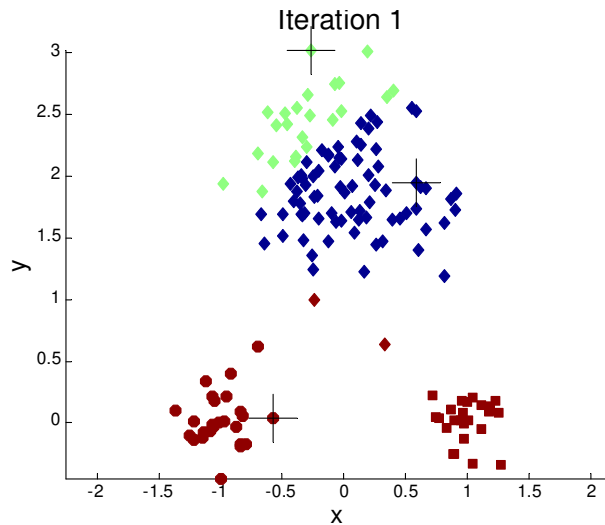


SUB-OPTIMAL CLUSTERING

# Importance of Choosing Initial Centroids



# Importance of Choosing Initial Centroids





# 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

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

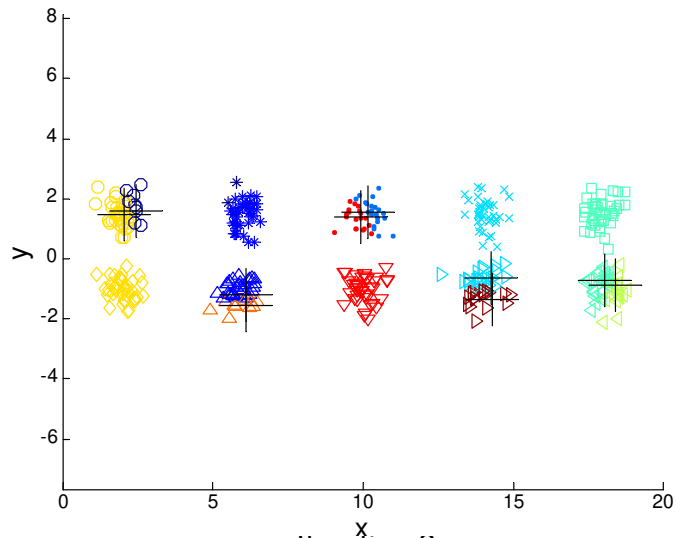
✗ For example, if  $K = 10$ , then probability =  $10!/10^{10} = 0.00036$

✗ Sometimes the initial centroids will readjust themselves in 'right' way, and sometimes they don't

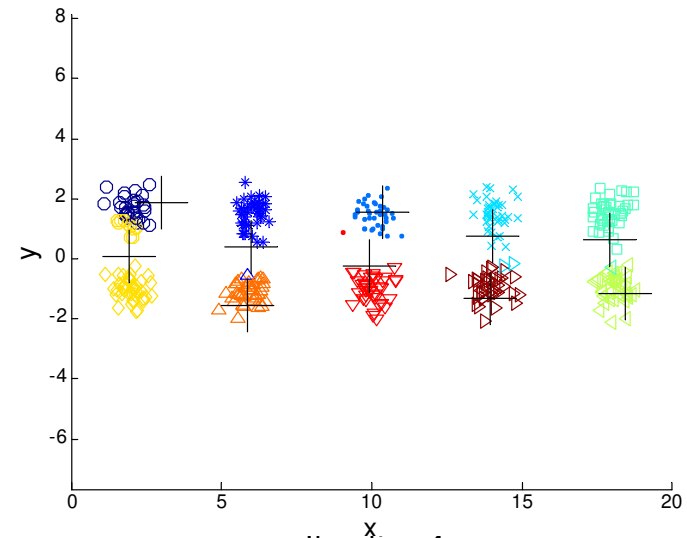
✗ Consider an example of five pairs of clusters

# 10 Clusters Example

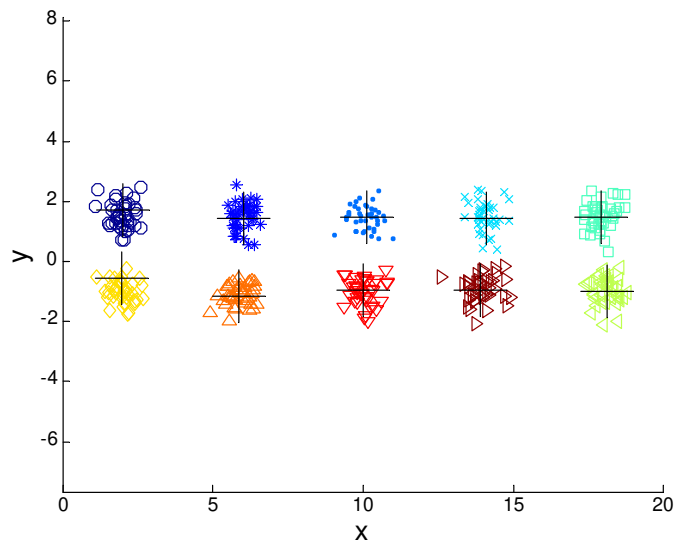
Iteration 1



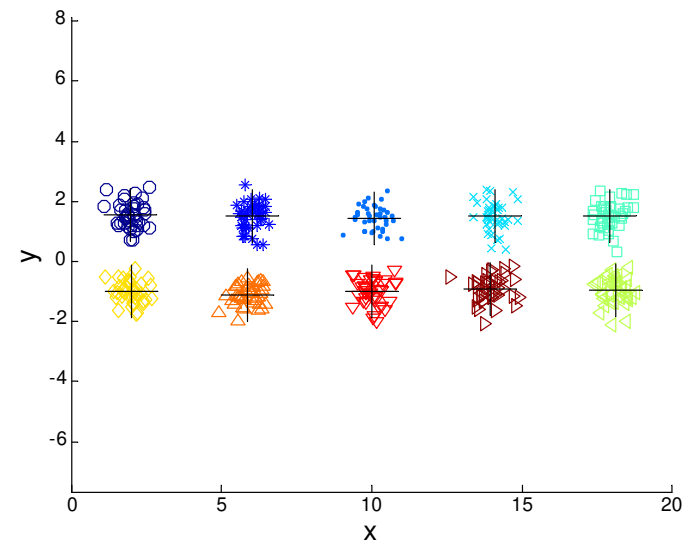
Iteration 2



Iteration 3

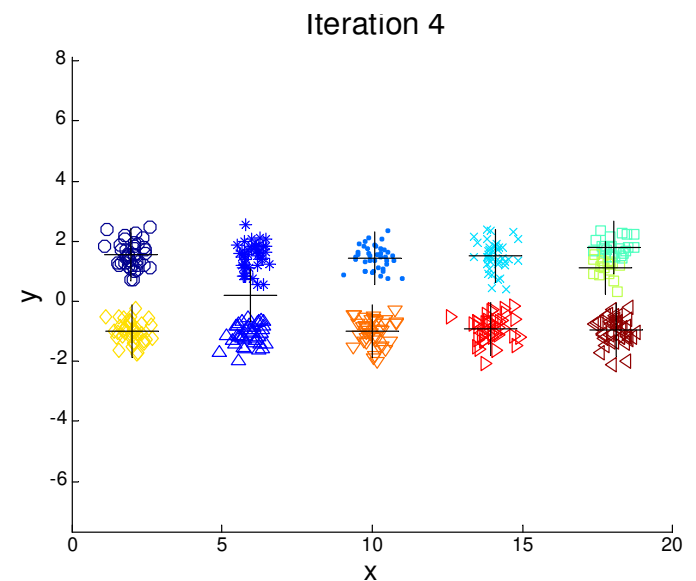
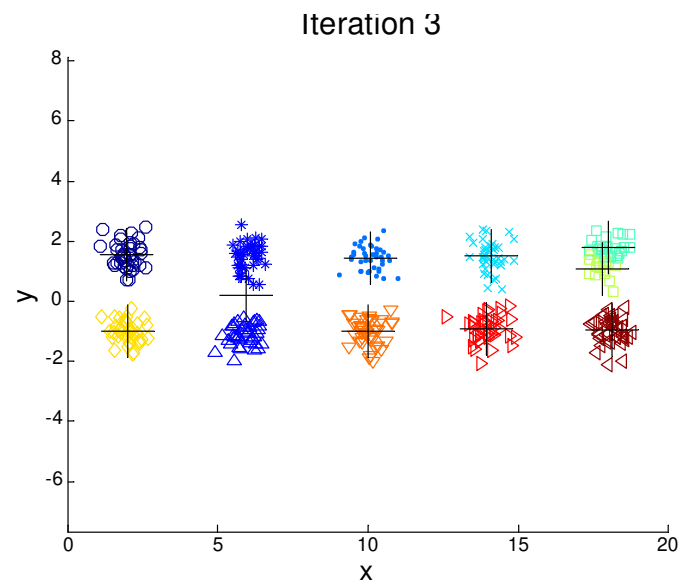
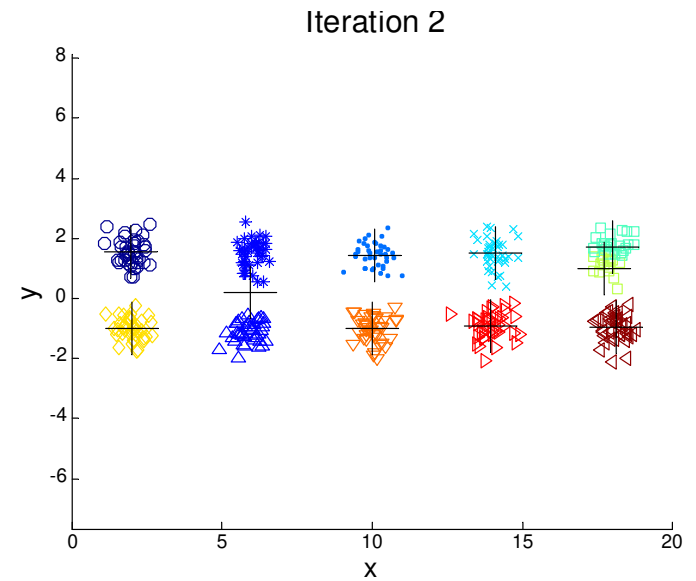
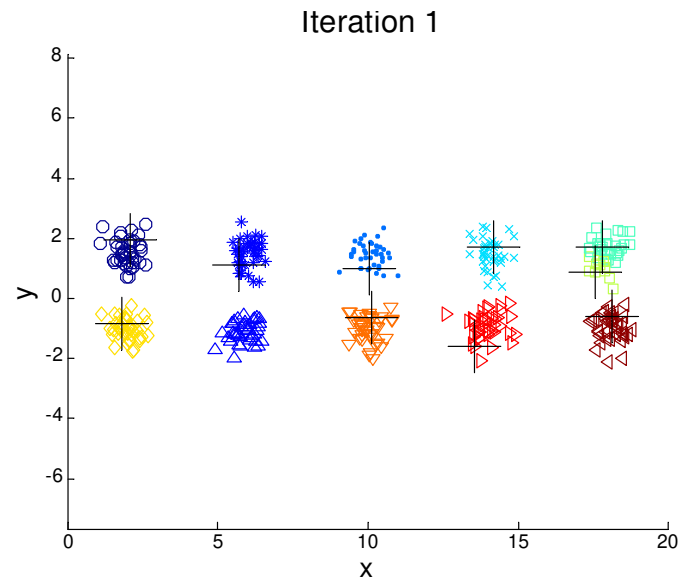


Iteration 4



STARTING WITH TWO INITIAL CENTROIDS IN ONE CLUSTER OF EACH PAIR OF CLUSTERS

# 10 Clusters Example



STARTING WITH SOME PAIRS OF CLUSTERS HAVING THREE INITIAL CENTROIDS, WHILE OTHER HAVE ONLY ONE.

# 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

- ✗ Split/merge clusters

- ✗ BISECTING K-MEANS

- ✗ Not as susceptible to initialization issues

# Handling Empty Clusters

- ✗ BASIC K-MEANS ALGORITHM CAN YIELD EMPTY CLUSTERS
- ✗ SEVERAL STRATEGIES
  - ✗ Choose the point that contributes most to SSE
  - ✗ Choose a point from the cluster with the highest SSE
  - ✗ If there are several empty clusters, the above can be repeated several times.

# Updating Centers Incrementally

- ✗ IN THE BASIC K-MEANS ALGORITHM, CENTROIDS ARE UPDATED AFTER ALL POINTS ARE ASSIGNED TO A CENTROID
- ✗ AN ALTERNATIVE IS TO UPDATE THE CENTROIDS AFTER EACH ASSIGNMENT (INCREMENTAL APPROACH)
  - ✗ Each assignment updates zero or two centroids
  - ✗ More expensive
  - ✗ Introduces an order dependency
  - ✗ Never get an empty cluster
  - ✗ Can use “weights” to change the impact

# 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
  - ✕ ISODATA (very **complex** and **maaaaaany** hyperparameters)

# Bisecting K-means

- ✕BISECTING K-MEANS ALGORITHM

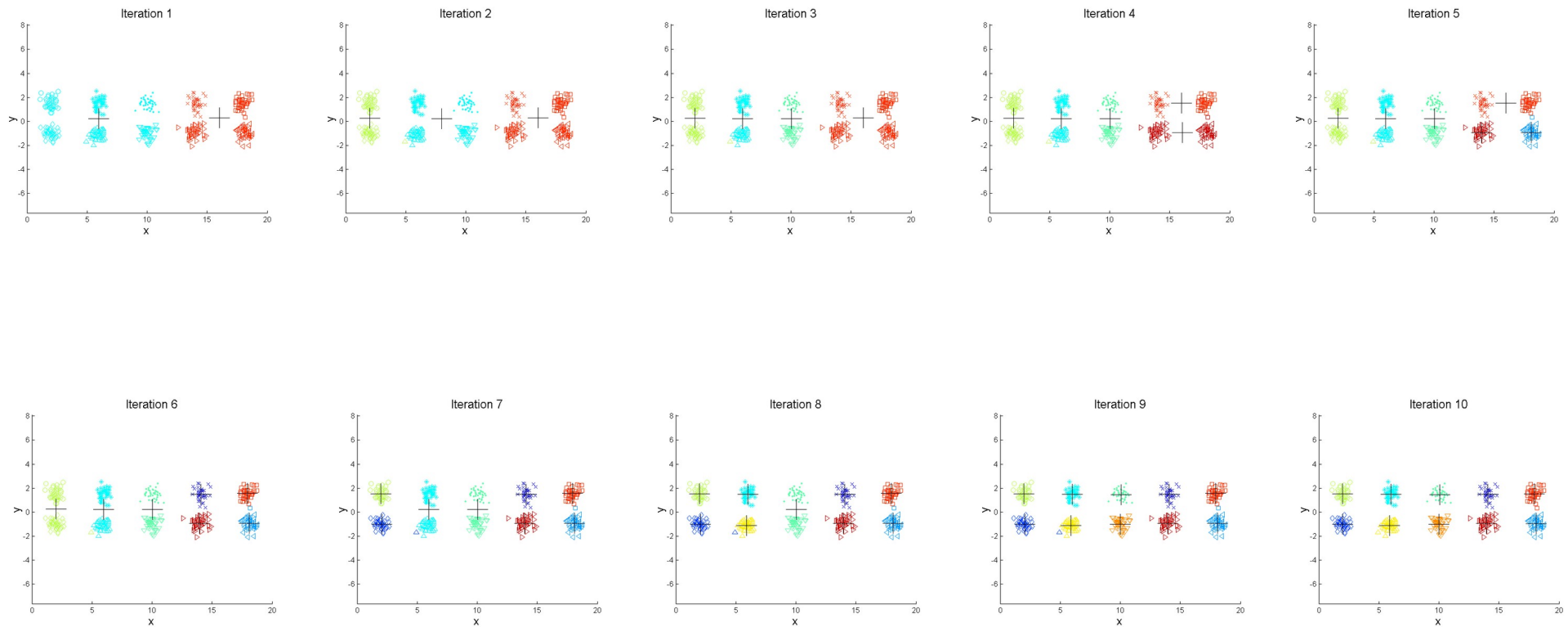
- ✕Variant of K-means that can produce a partitional or a hierarchical clustering



# Bisecting k-means algorithm

- 1) START WITH A CLUSTER WITH ALL THE DATA POINTS
- 2) REPEAT
  - 2.1) Pick a cluster to split: Choose an appropriate measure of sparseness
  - 2.2) Find 2-sub clusters using the k-means algorithm
  - 2.3) Repeat step 2.2 for ITER times and take the split that produces the clustering with the highest overall similarity
- 3) UNTIL THE DESIRED NUMBER OF CLUSTERS IS REACHED

# Bisecting K-means Example



# 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.

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

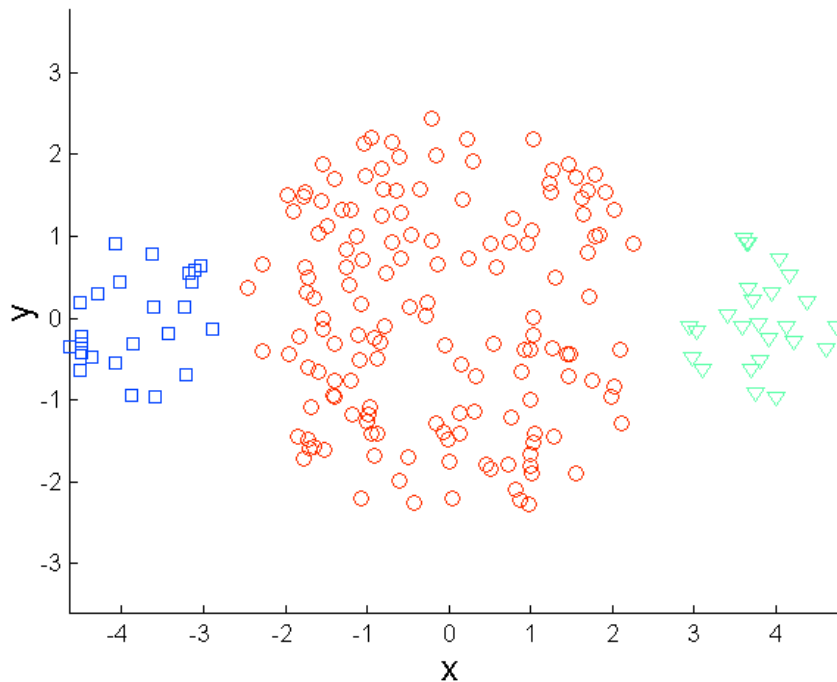
- ✕  $x$  is a data point in cluster  $C_i$  and  $m_i$  is the representative point for cluster  $C_i$ 
  - ✕ can show that  $m_i$  corresponds to the center (mean) of the cluster
- ✕ Given two clusters, we can choose the one with the smallest error
- ✕ One easy way to reduce SSE is to increase  $K$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

# Limitations of K-means

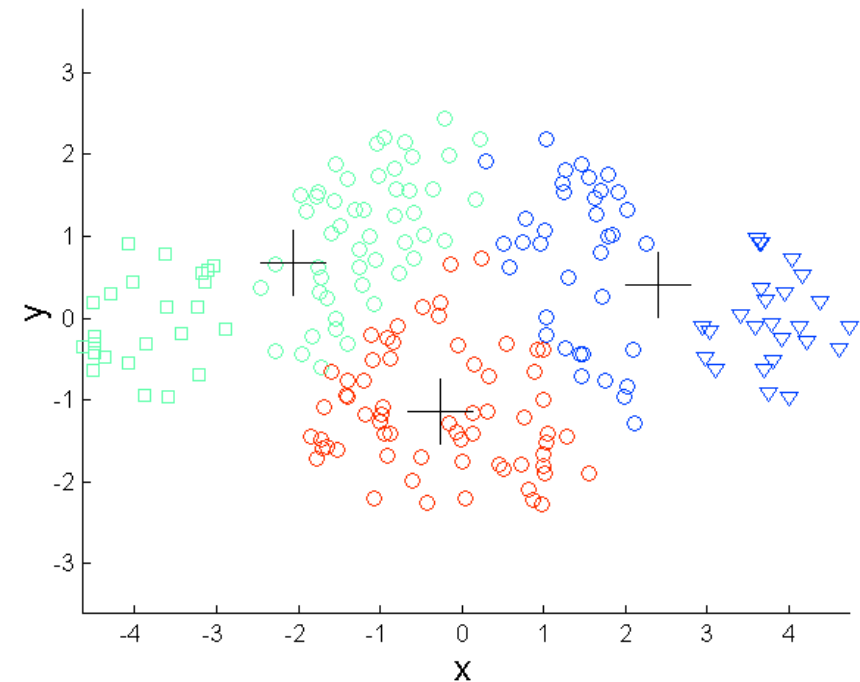
- ✗ K-MEANS HAS PROBLEMS WHEN CLUSTERS ARE OF DIFFERING
  - ✗ Sizes
  - ✗ Densities
  - ✗ Non-globular shapes
- ✗ K-MEANS HAS PROBLEMS WHEN THE DATA CONTAINS OUTLIERS

# Limitations of K-means: Differing Sizes

**X**TENDENCY TO BREAK LARGE CLUSTERS



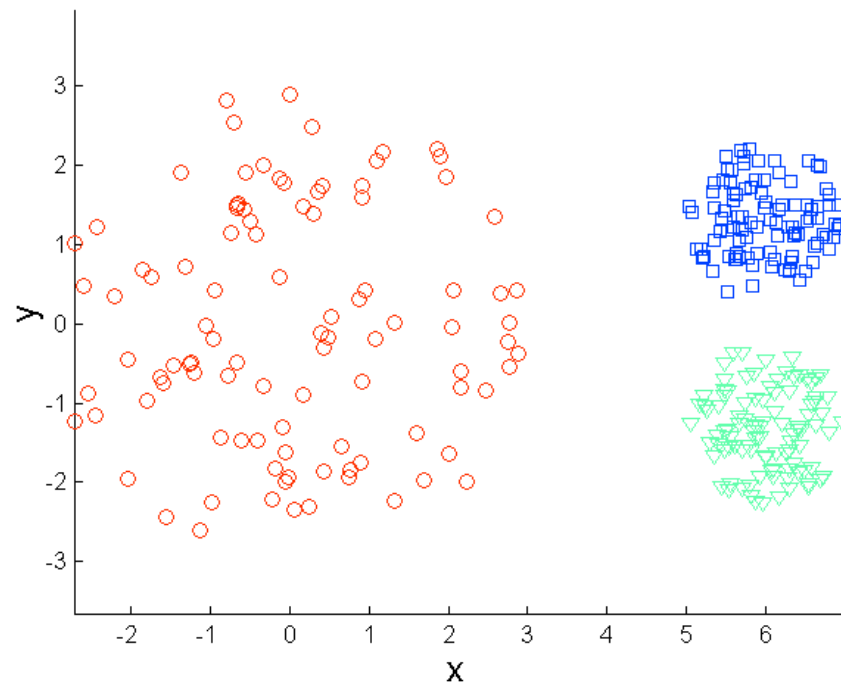
ORIGINAL POINTS



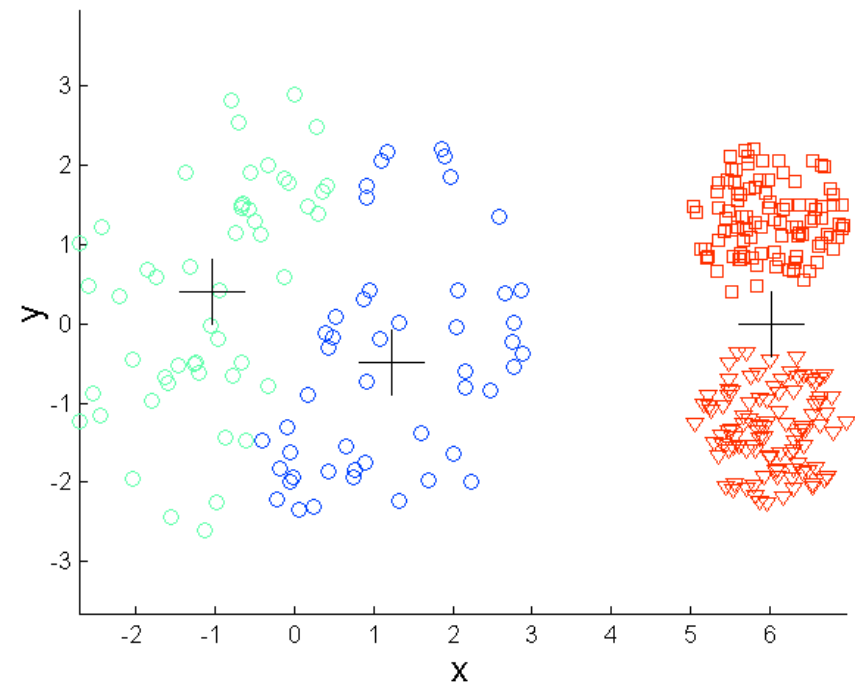
K-MEANS (3 CLUSTERS)

# Limitations of K-means: Differing Density

✗ DENSITY CONCEPT IS NOT CONSIDERED BY K-MEANS

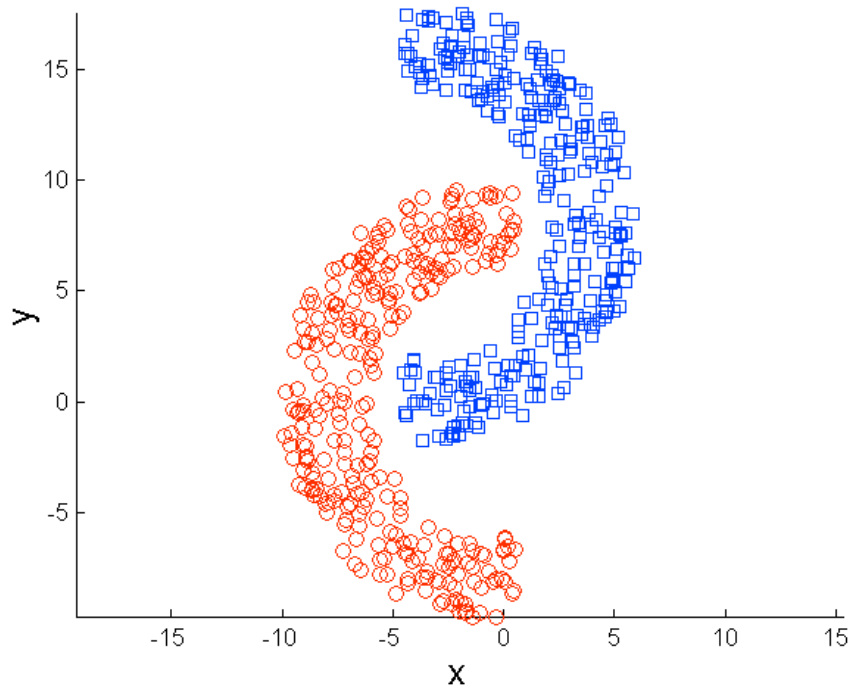


ORIGINAL POINTS

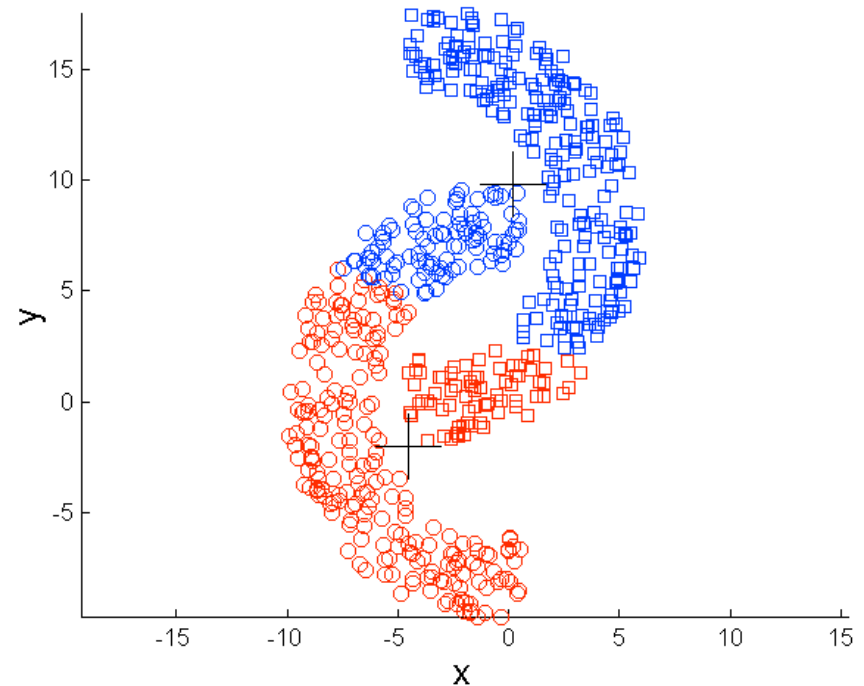


K-MEANS (3 CLUSTERS)

# Limitations of K-means: Non-globular Shapes

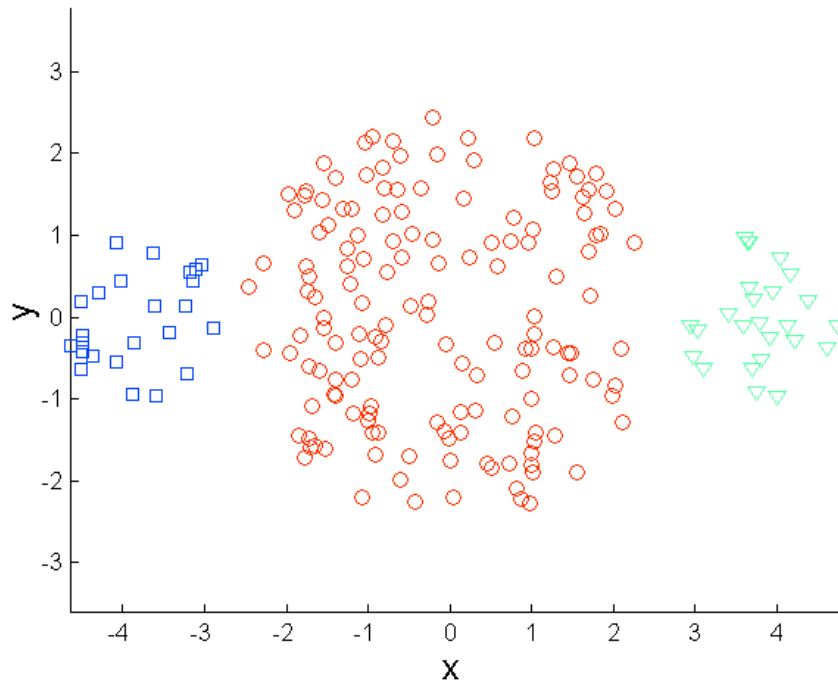


ORIGINAL POINTS

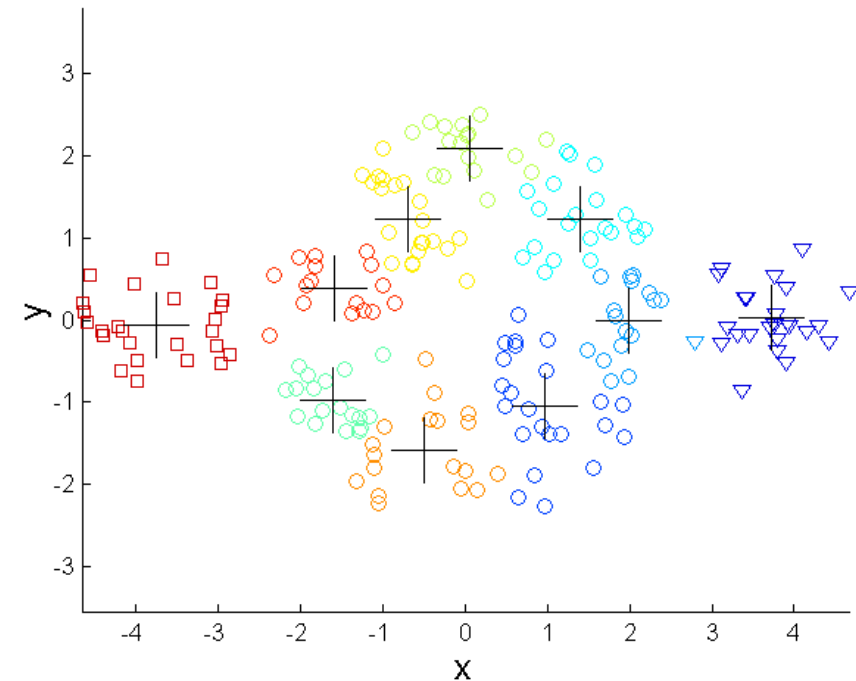


K-MEANS (2 CLUSTERS)

# Overcoming K-means Limitations



ORIGINAL POINTS



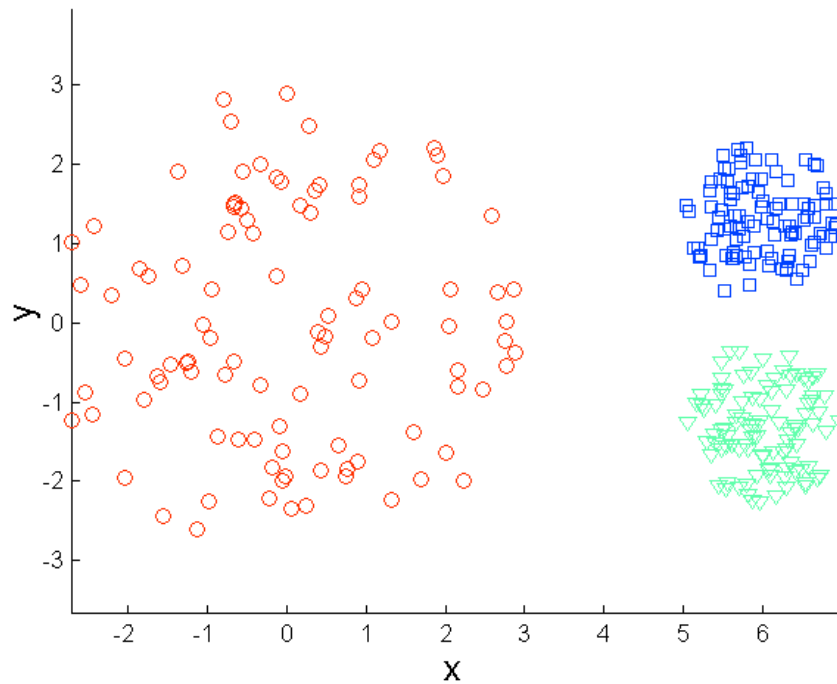
K-MEANS CLUSTERS

ONE SOLUTION IS TO USE MANY CLUSTERS.

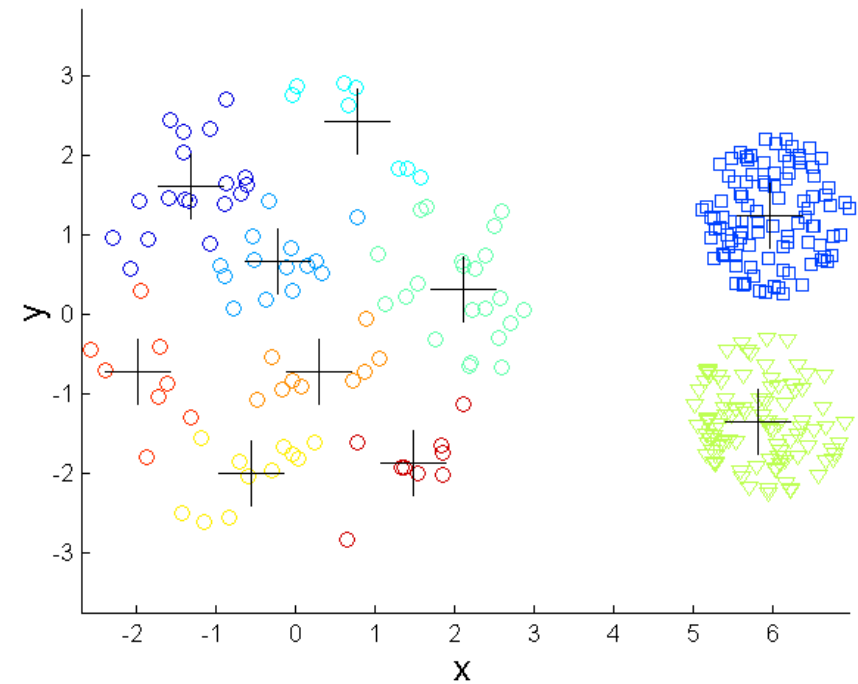
FIND PARTS OF CLUSTERS, BUT NEED TO PUT TOGETHER.



# Overcoming K-means Limitations

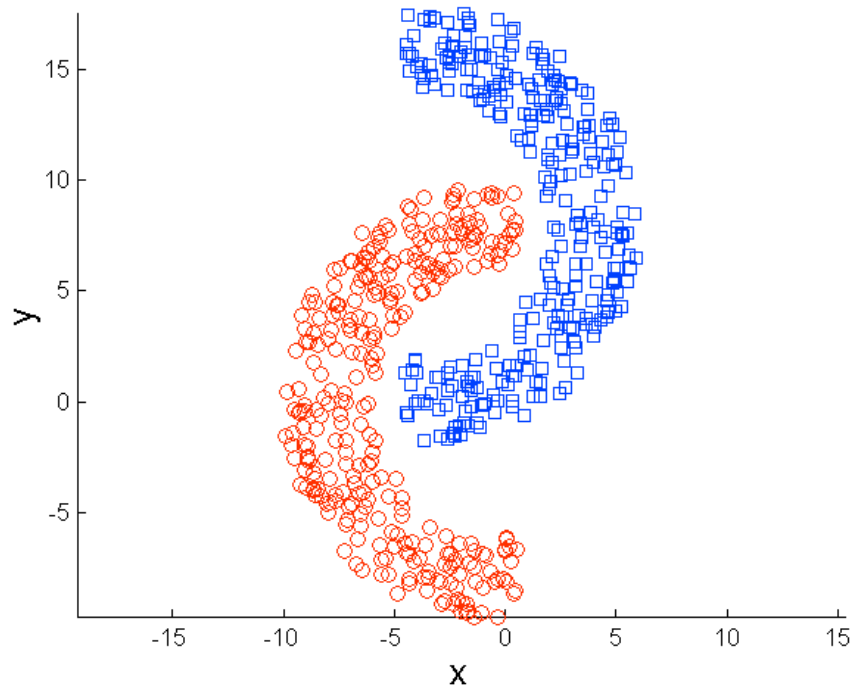


ORIGINAL POINTS

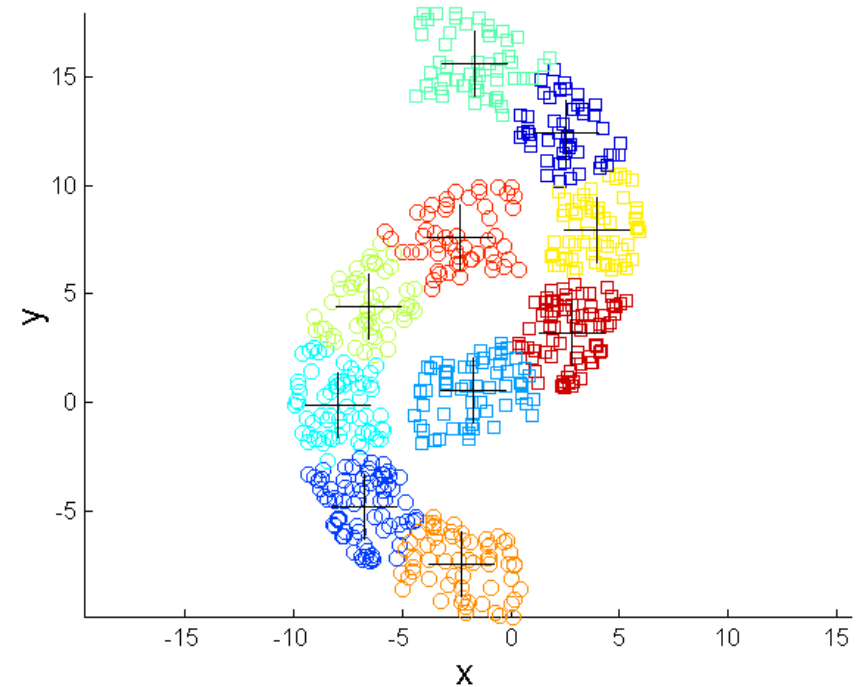


K-MEANS CLUSTERS

# Overcoming K-means Limitations



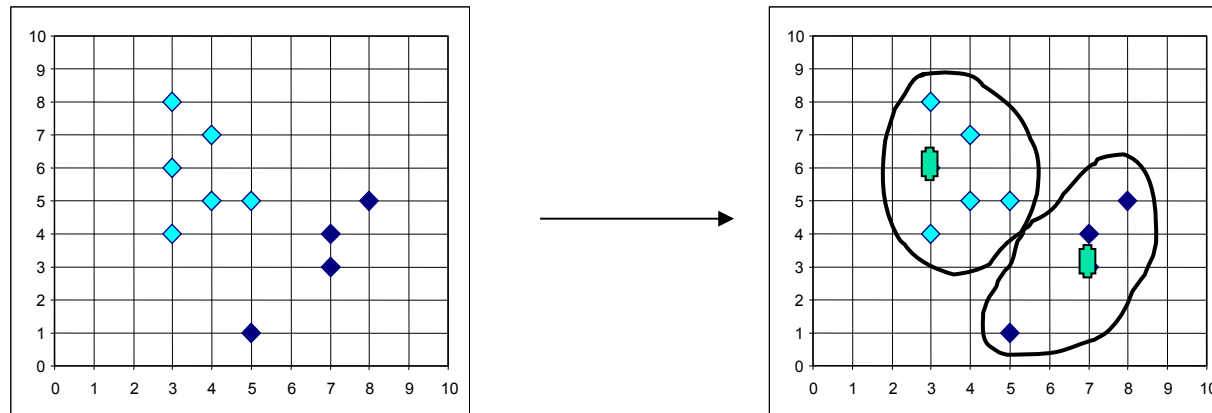
ORIGINAL POINTS



K-MEANS CLUSTERS

# k-Means and Outliers

- ✗ THE K-MEANS ALGORITHM IS SENSITIVE TO OUTLIERS !
  - ✗ Since an object with an extremely large value may substantially distort the distribution of the data.
- ✗ K-MEDOIDS: INSTEAD OF TAKING THE **MEAN** VALUE OF THE OBJECT IN A CLUSTER AS A REFERENCE POINT, **MEDOIDS** CAN BE USED, WHICH IS THE **MOST CENTRALLY LOCATED** OBJECT IN A CLUSTER.



# The *K-Medoids* Clustering Method

- ✕ FIND REPRESENTATIVE OBJECTS, CALLED MEDOIDS, IN CLUSTERS
- ✕ PAM (PARTITIONING AROUND MEDOIDS, 1987)
  - ✕ starts from an initial set of medoids and iteratively replaces one of the medoids by one of the non-medoids if it improves the total distance of the resulting clustering
  - ✕ PAM works effectively for small data sets, but does not scale well for large data sets
- ✕ CLARA (KAUFMANN & ROUSSEEUW, 1990)
- ✕ CLARANS (NG & HAN, 1994): RANDOMIZED SAMPLING
- ✕ FOCUSING + SPATIAL DATA STRUCTURE (ESTER ET AL., 1995)

# PAM Algorithm

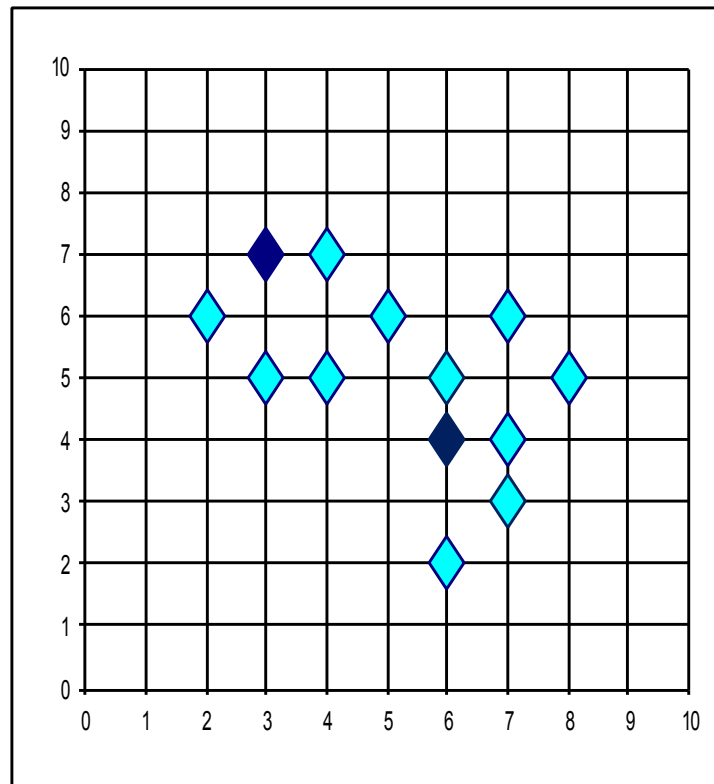
- ✕ USE REAL OBJECT TO REPRESENT THE CLUSTER
  - ✕ Select  $k$  representative objects arbitrarily
  - ✕ For each pair of non-selected object  $h$  and selected object  $i$ , calculate the total swapping cost  $TC_{ih} = \sum_j C_{jih}$
  - ✕ For each pair of  $i$  and  $h$ ,
    - ✕ If  $TC_{ih} < 0$ ,  $i$  is replaced by  $h$
    - ✕ Then assign each non-selected object to the most similar representative object
  - ✕ repeat steps 2-3 until there is no change

$C_{jih}$  is the cost of swapping  $i$  with  $h$  for all non medoid objects  $j$

# Case (i) : Computation of $C_{jih}$

- $j$  currently belongs to the cluster represent by the medoid  $i$
- $j$  is less similar to the medoid  $t$  compare to  $h$

Therefore, in future,  $j$  belongs to the cluster represented by  $h$



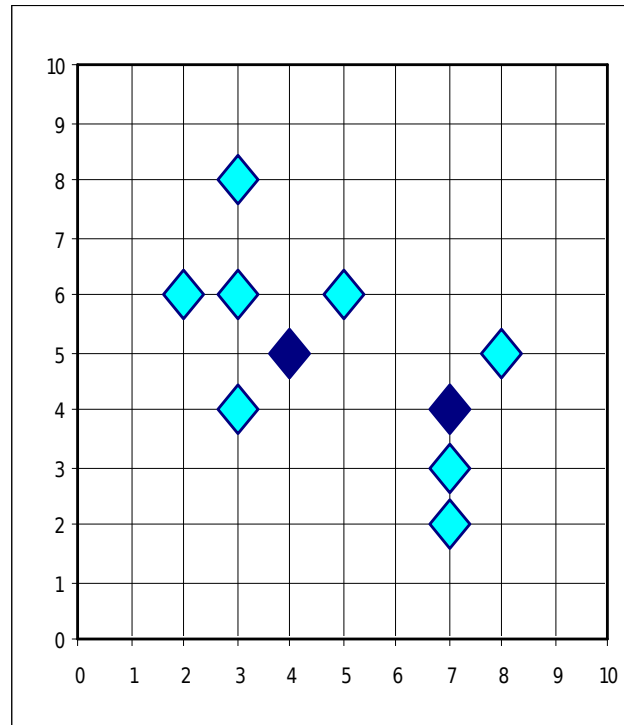
$$E = \sum_{i=1}^k \sum_{j \in C_i} d(j, i)$$

$$C_{jih} = d(j, h) - d(j, i)$$

# Case (ii) : Computation of $C_{jih}$

- $j$  currently belongs to the cluster represent by the medoid  $i$
- $j$  is more similar to the medoid  $t$  compare to  $h$

Therefore, in future,  $j$  belongs to the cluster represented by  $t$

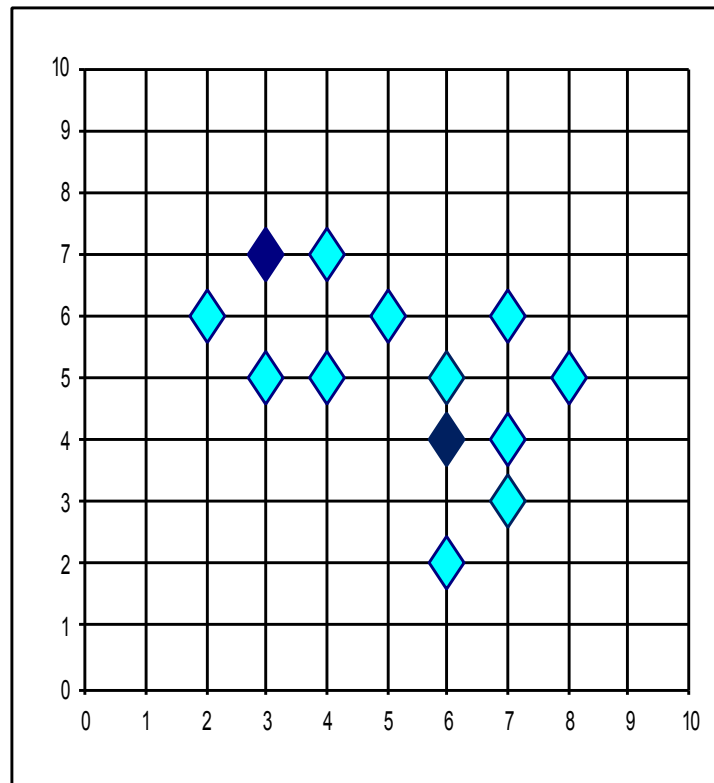


$$C_{jih} = d(j, t) - d(j, i)$$

# Case (iii) : Computation of $C_{jih}$

- $j$  currently belongs to the cluster represent by the medoid  $t$
- $j$  is more similar to the medoid  $t$  compare to  $h$

Therefore, in future,  $j$  belongs to the cluster represented by  $t$  itself



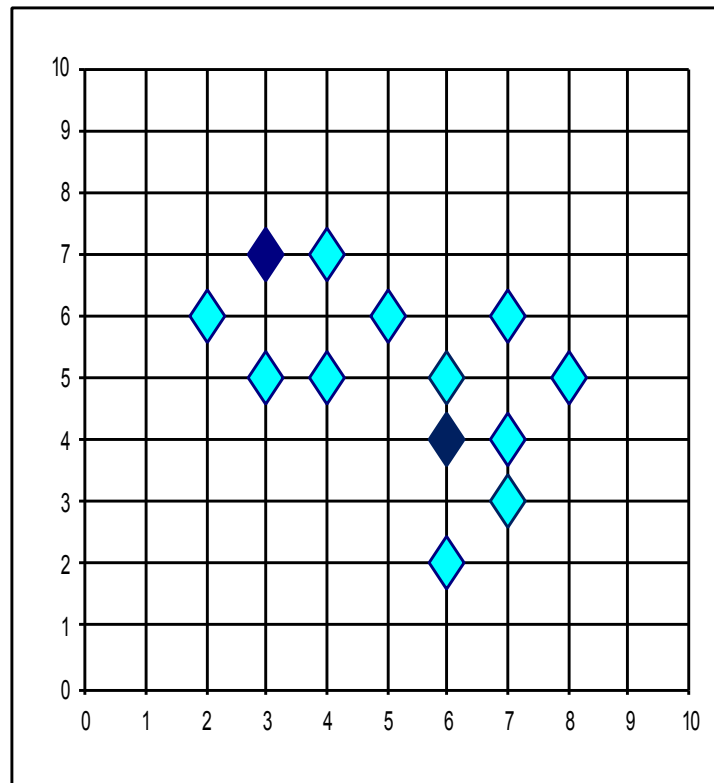
$$C_{jih} = 0$$



# Case (iv) : Computation of $C_{jih}$

- $j$  currently belongs to the cluster represent by the medoid  $t$
- $j$  is less similar to the medoid  $t$  compare to  $h$

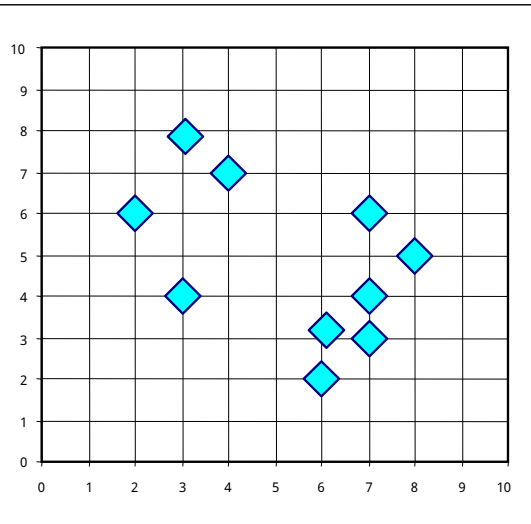
Therefore, in future,  $j$  belongs to the cluster represented by  $h$



$$C_{jih} = d(j, t) - d(j, h)$$

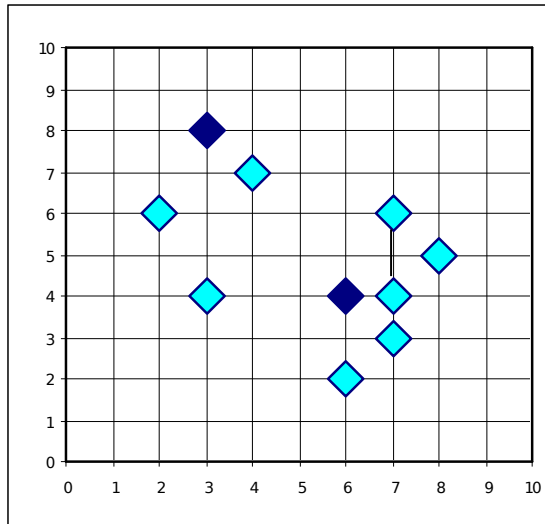
# A Typical K-Medoids Algorithm (PAM)

Total Cost = 20

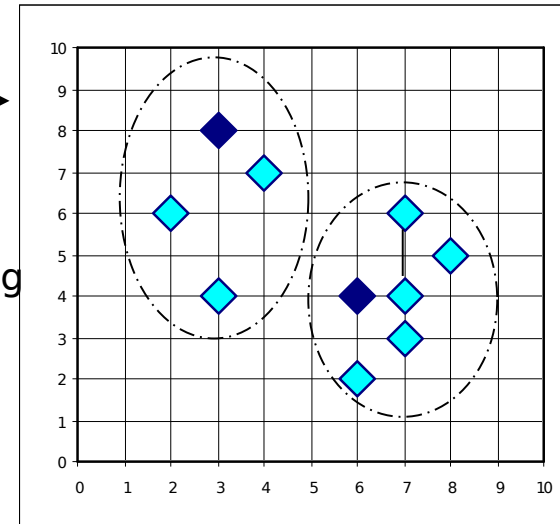


K=2

Arbitrary  
choose k  
object as  
initial  
medoids



Assign  
each  
remaining  
object to  
nearest  
medoids

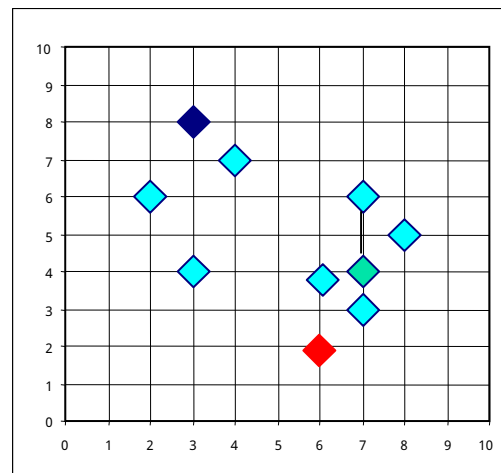


Randomly select a  
nonmedoid object,  $O_{\text{random}}$

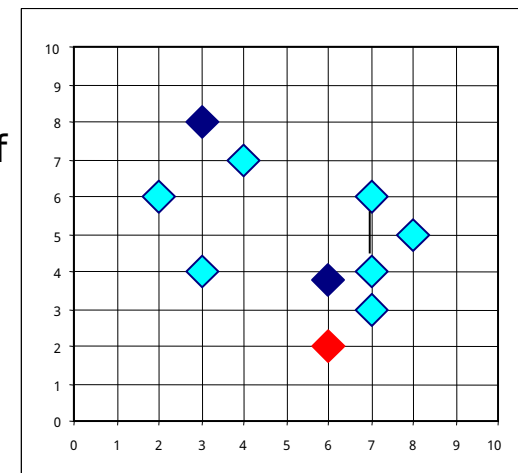
Total Cost = 26

**Do loop  
Until no  
change**

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



Compute  
total cost of  
swapping



# What Is the Problem with PAM?

- ✗ PAM IS MORE ROBUST THAN K-MEANS IN THE PRESENCE OF NOISE AND OUTLIERS BECAUSE A MEDOID IS LESS INFLUENCED BY OUTLIERS OR OTHER EXTREME VALUES THAN A MEAN
- ✗ PAM WORKS EFFICIENTLY FOR SMALL DATA SETS BUT DOES NOT **SCALE WELL** FOR LARGE DATA SETS.

✗  $O(k(n-k)^2)$  for each iteration

WHERE N IS # OF DATA, K IS # OF CLUSTERS

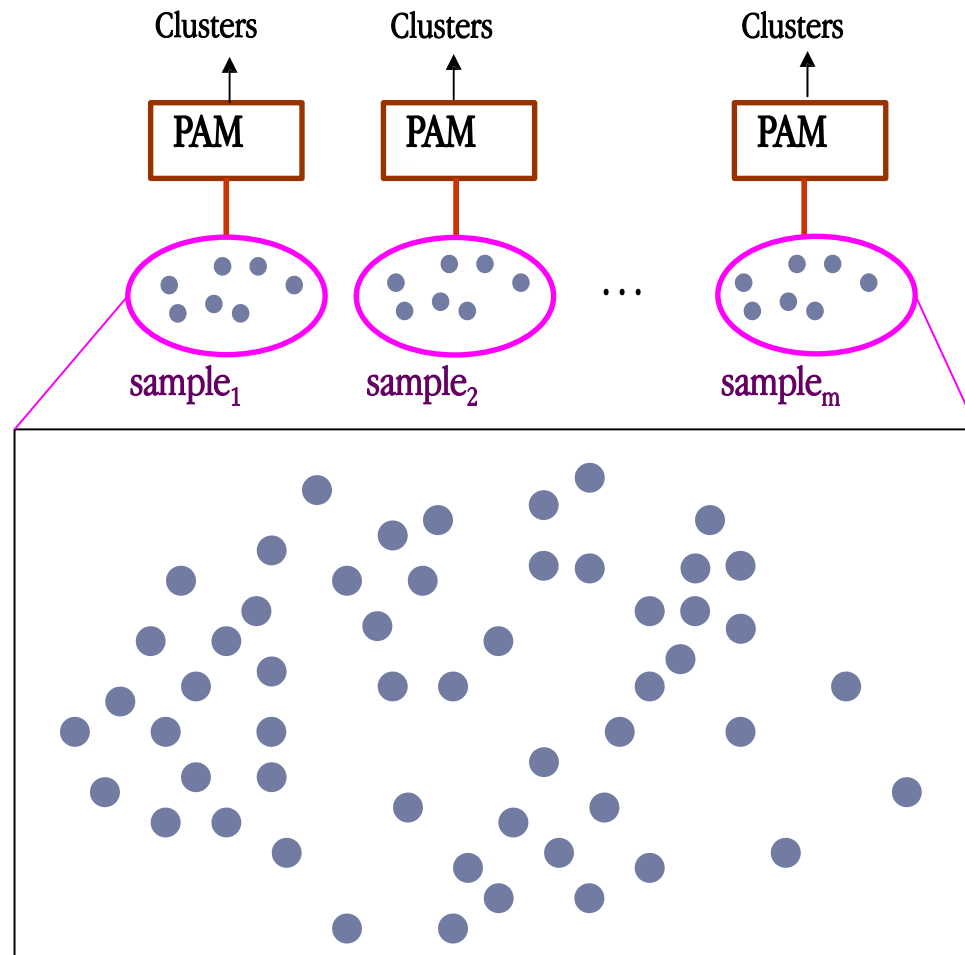
- ➔ SAMPLING BASED METHOD,  
CLARA(CLUSTERING LARGE APPLICATIONS)

# CLARA (Clustering Large Applications)

- ✕ CLARA (KAUFMANN AND ROUSSEEUW IN 1990)
- ✕ IT DRAWS MULTIPLE SAMPLES OF THE DATA SET, APPLIES *PAM* ON EACH SAMPLE, AND GIVES THE BEST CLUSTERING AS THE OUTPUT
- ✕ STRENGTH: DEALS WITH LARGER DATA SETS THAN *PAM*
- ✕ WEAKNESS:
  - ✕ Efficiency depends on the sample size
  - ✕ A good clustering based on samples will not necessarily represent a good clustering of the whole data set if the sample is biased

# CLARA

Choose the best clustering



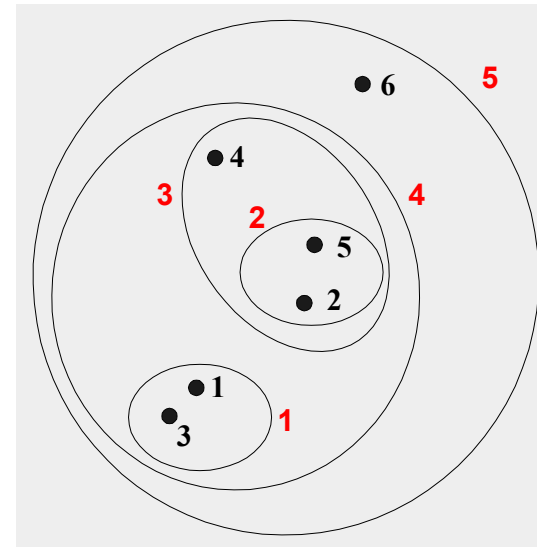
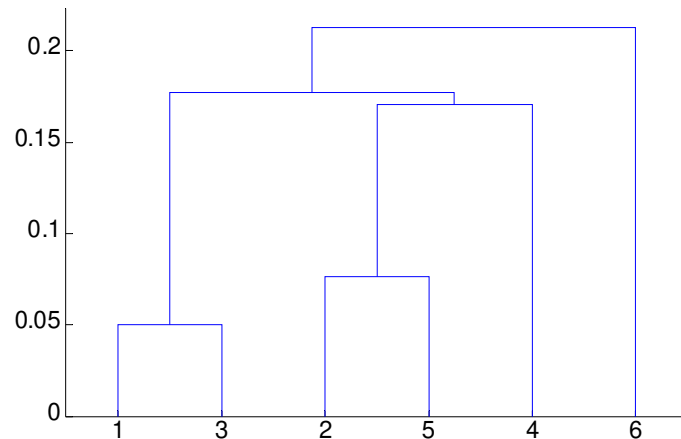
# *CLARANS* (“Randomized” CLARA) (1994)

- ✕ *CLARANS* (A CLUSTERING ALGORITHM BASED ON RANDOMIZED SEARCH) (NG AND HAN'94)
- ✕ *CLARANS* DRAWS SAMPLE OF NEIGHBORS DYNAMICALLY
- ✕ THE CLUSTERING PROCESS CAN BE PRESENTED AS SEARCHING A GRAPH WHERE EVERY NODE IS A POTENTIAL SOLUTION, THAT IS, A SET OF  $K$  MEDOIDS
- ✕ IF THE LOCAL OPTIMUM IS FOUND, *CLARANS* STARTS WITH NEW RANDOMLY SELECTED NODE IN SEARCH FOR A NEW LOCAL OPTIMUM
- ✕ IT IS MORE EFFICIENT AND SCALABLE THAN BOTH *PAM* AND *CLARA*
- ✕ FOCUSING TECHNIQUES AND SPATIAL ACCESS STRUCTURES MAY FURTHER IMPROVE ITS PERFORMANCE (ESTER ET AL.'95)

# Hierarchical clustering

# 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 MATRIX

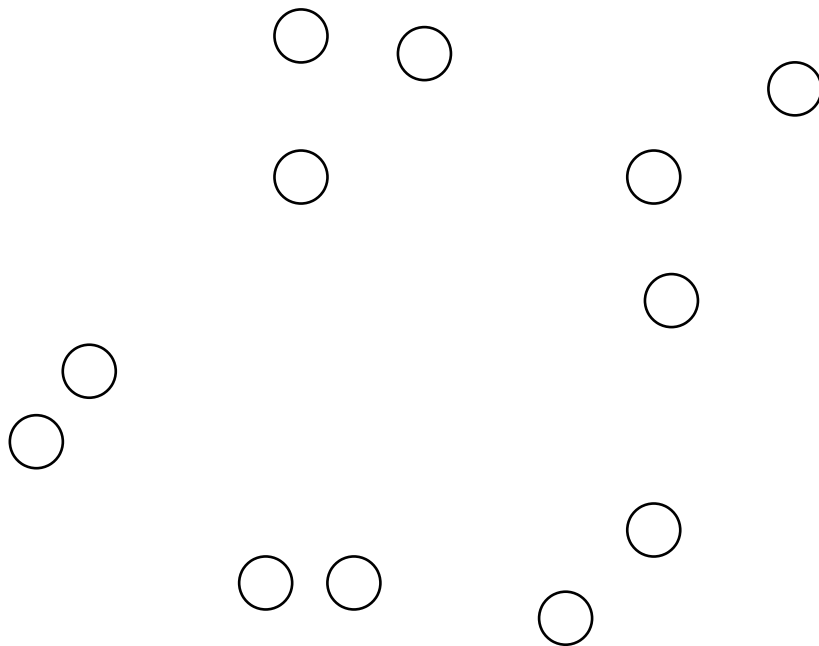
- ✕Merge or split one cluster at a time

# Agglomerative Clustering Algorithm

- ✗ MORE POPULAR HIERARCHICAL CLUSTERING TECHNIQUE
- ✗ BASIC ALGORITHM IS STRAIGHTFORWARD
  1. Compute the proximity matrix
  2. Let each data point be a cluster
  3. **Repeat**
    4. Merge the two closest clusters
    5. Update the proximity matrix
  6. **Until** only a single cluster remains
- ✗ KEY OPERATION IS THE COMPUTATION OF THE PROXIMITY OF TWO CLUSTERS
  - ✗ Different approaches to defining the distance between clusters distinguish the different algorithms

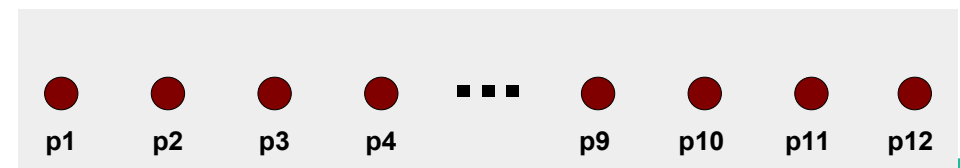
# Starting Situation

✕ START WITH CLUSTERS OF INDIVIDUAL POINTS AND A PROXIMITY MATRIX



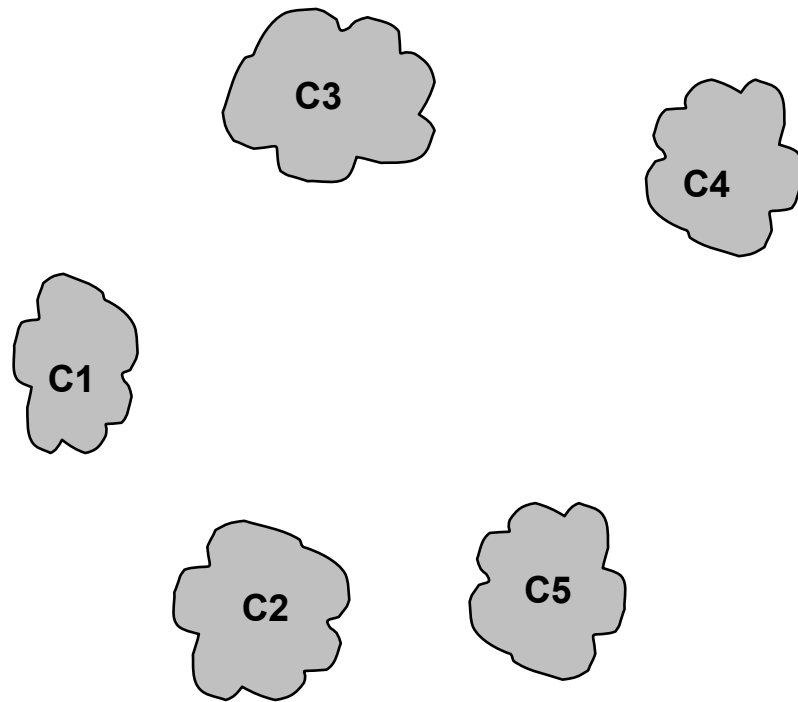
	p1	p2	p3	p4	p5	...
p1						
p2						
p3						
p4						
p5						
.						
.						
.						

PROXIMITY MATRIX



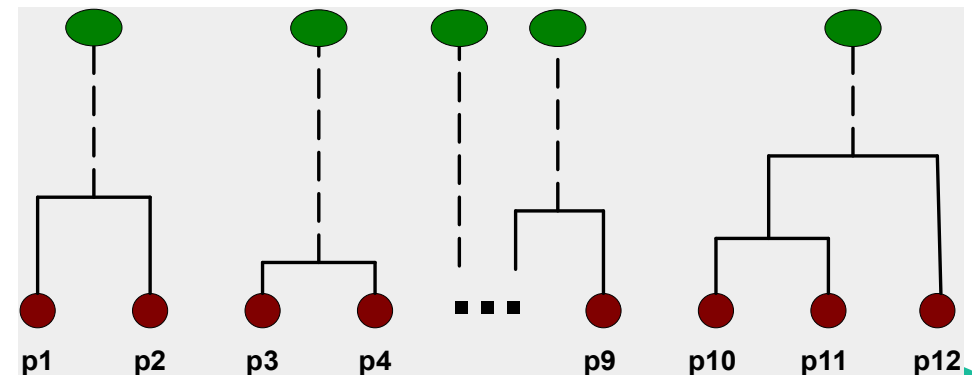
# Intermediate Situation

✗ AFTER SOME MERGING STEPS, WE HAVE SOME CLUSTERS



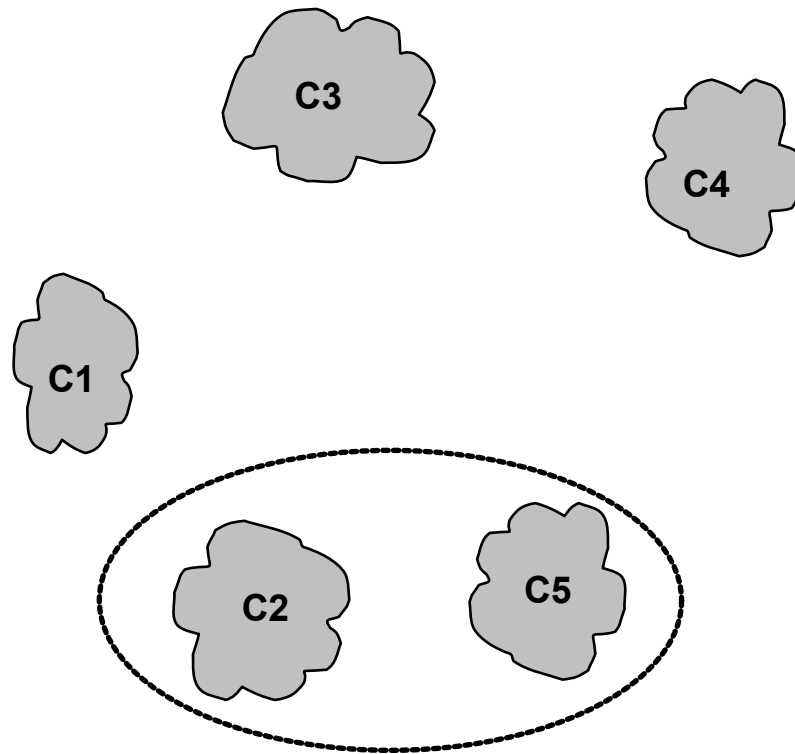
	C1	C2	C3	C4	C5
C1					
C2					
C3					
C4					
C5					

PROXIMITY MATRIX



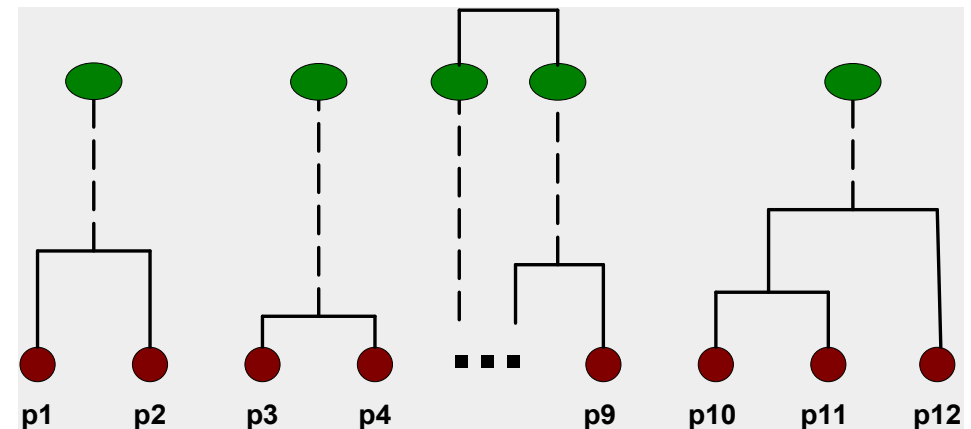
# Intermediate Situation

✗ WE WANT TO MERGE THE TWO CLOSEST CLUSTERS (C2 AND C5) AND UPDATE THE PROXIMITY MATRIX.



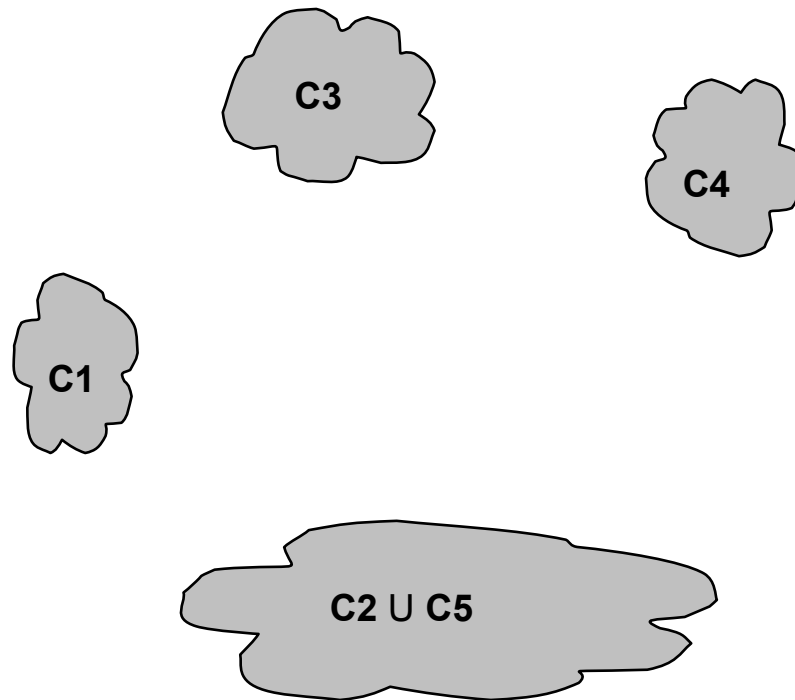
	C1	C2	C3	C4	C5
C1					
C2					
C3					
C4					
C5					

PROXIMITY MATRIX



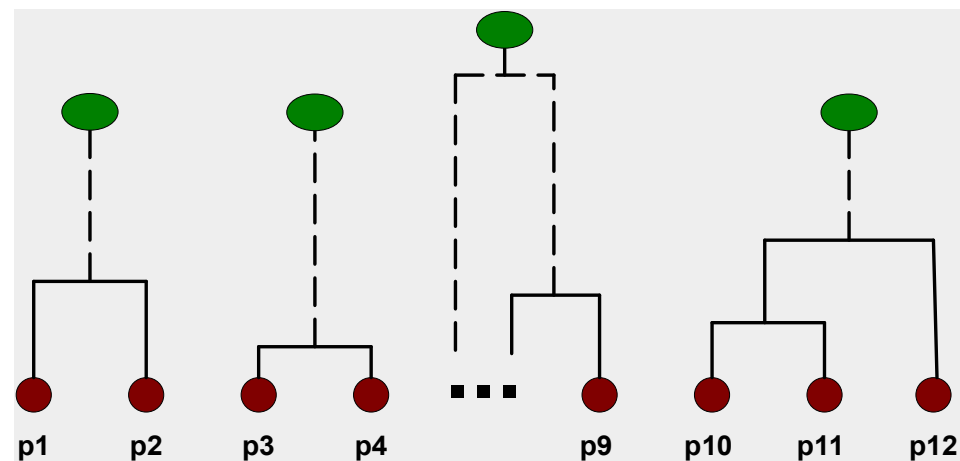
# After Merging

THE QUESTION IS “HOW DO WE UPDATE THE PROXIMITY MATRIX?”

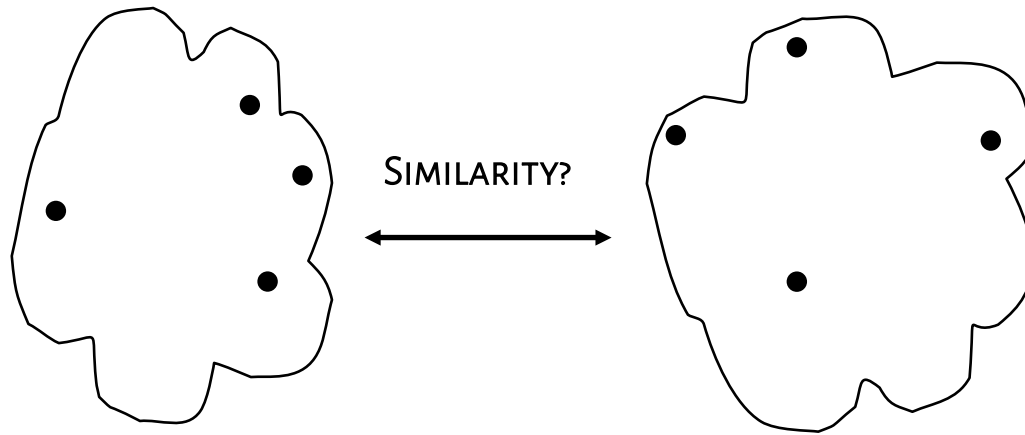


		C2 U C5			
		C1		C3	C4
C2 U C5	C1		?		
	C3		?		
	C4		?		

PROXIMITY MATRIX



# How to Define Inter-Cluster Similarity



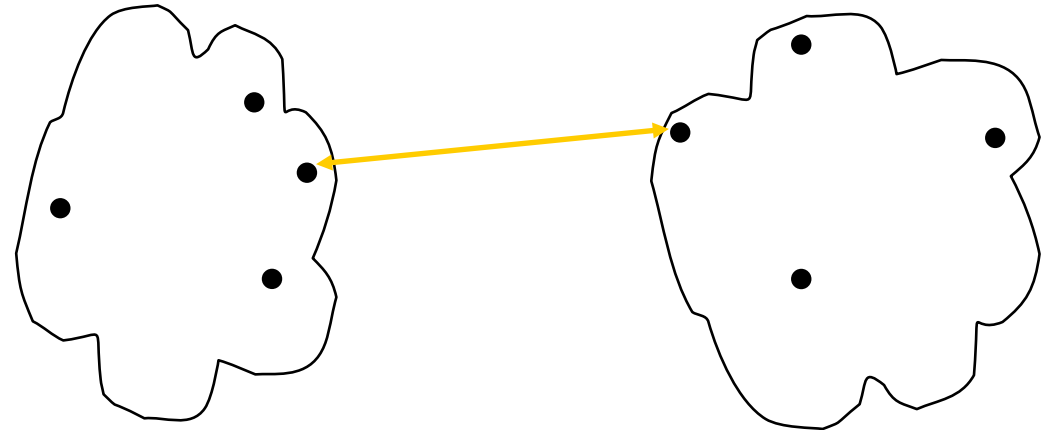
- MIN
- MAX
- Group Average
- Distance Between Centroids
- Other methods driven by an objective function
  - Ward's Method uses squared error

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

PROXIMITY MATRIX

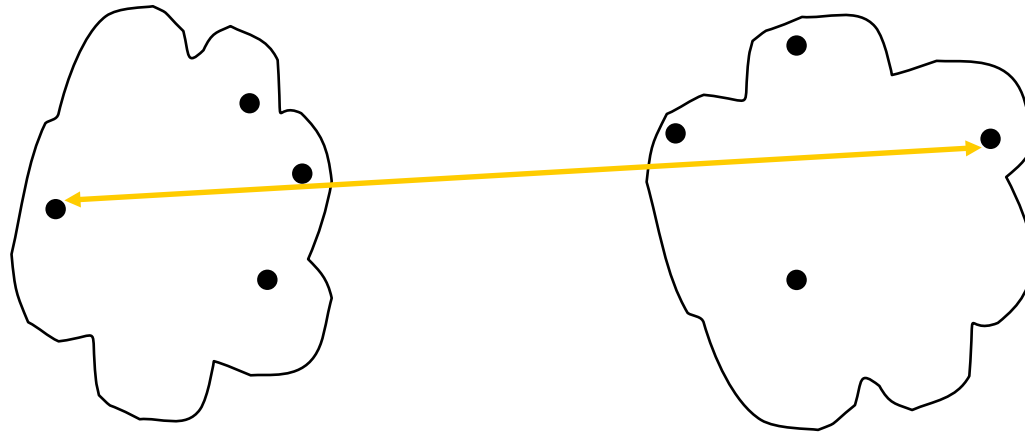


# How to Define Inter-Cluster Similarity



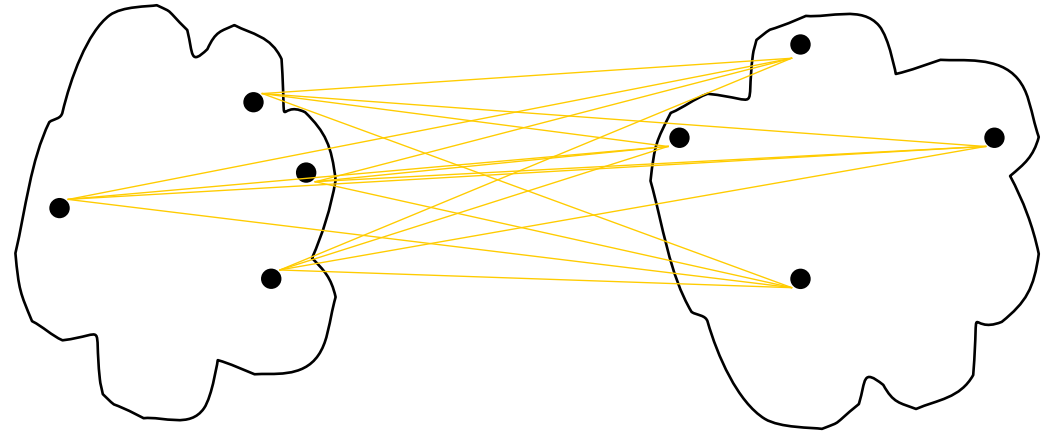
- MIN
- MAX
- Group Average
- Distance Between Centroids
- Other methods driven by an objective function
  - Ward's Method uses squared error

# How to Define Inter-Cluster Similarity



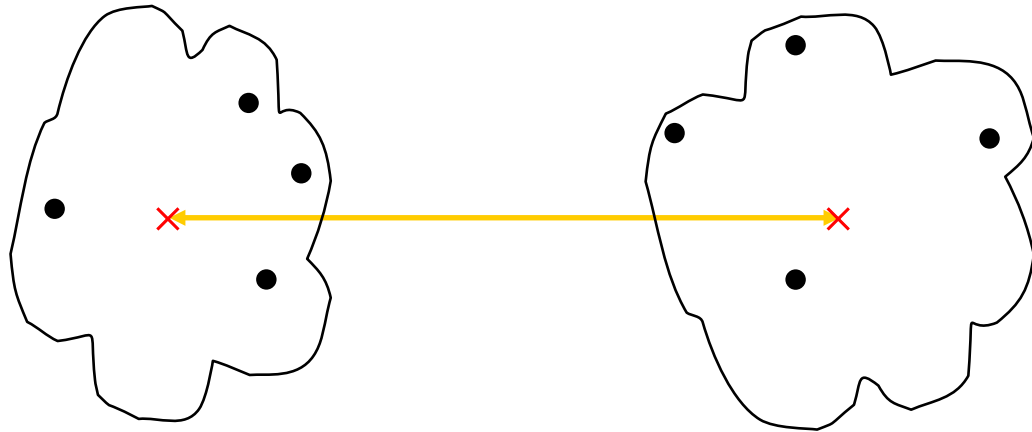
- MIN
- MAX
- Group Average
- Distance Between Centroids
- Other methods driven by an objective function
  - Ward's Method uses squared error

# How to Define Inter-Cluster Similarity



- MIN
- MAX
- Group Average
- Distance Between Centroids
- Other methods driven by an objective function
  - Ward's Method uses squared error

# How to Define Inter-Cluster Similarity

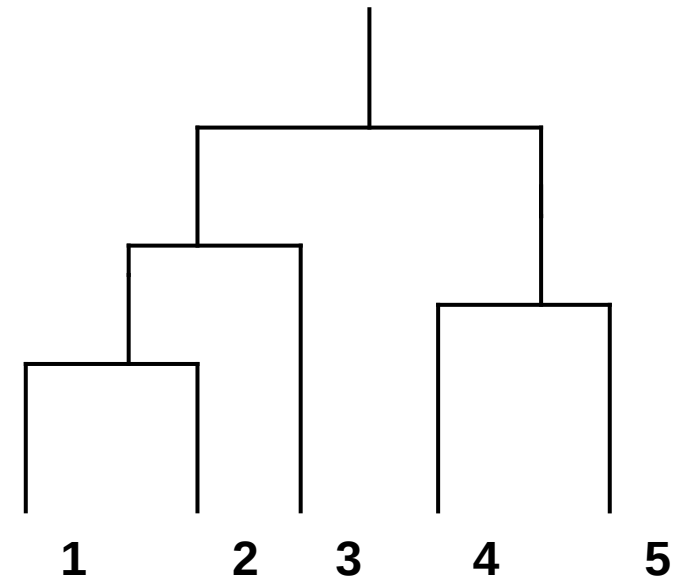


- MIN
- MAX
- Group Average
- Distance Between Centroids
- Other methods driven by an objective function
  - Ward's Method uses squared error

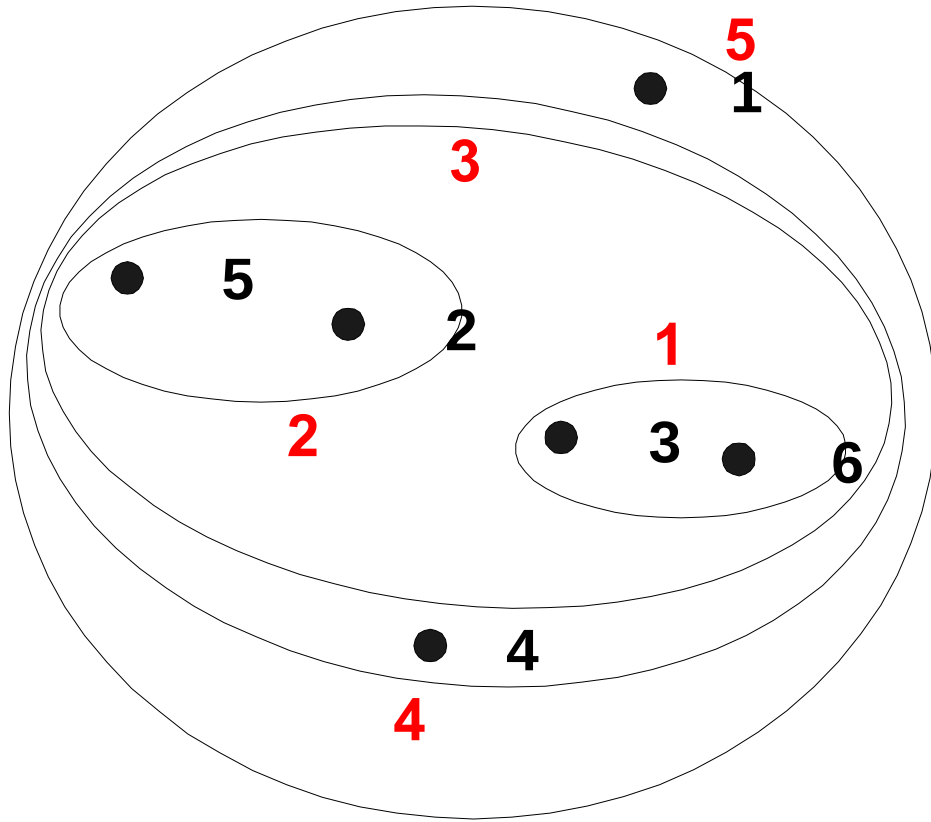
# Cluster Similarity: MIN or Single Link

- ✗ SIMILARITY OF TWO CLUSTERS IS BASED ON THE TWO MOST SIMILAR (CLOSEST) POINTS IN THE DIFFERENT CLUSTERS
  - ✗ Determined by one pair of points, i.e., by one link in the proximity graph.

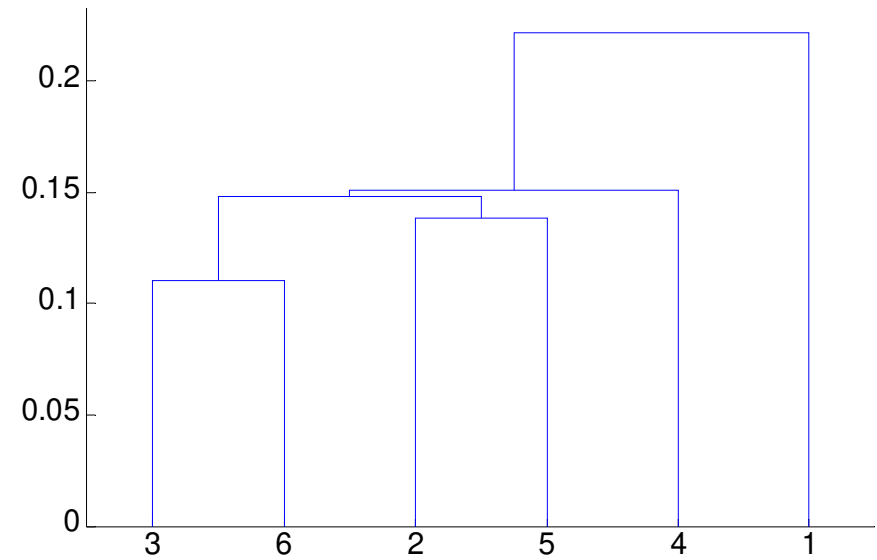
	I1	I2	I3	I4	I5
I1	1.00	0.90	0.10	0.65	0.20
I2	0.90	1.00	0.70	0.60	0.50
I3	0.10	0.70	1.00	0.40	0.30
I4	0.65	0.60	0.40	1.00	0.80
I5	0.20	0.50	0.30	0.80	1.00



# Hierarchical Clustering: MIN

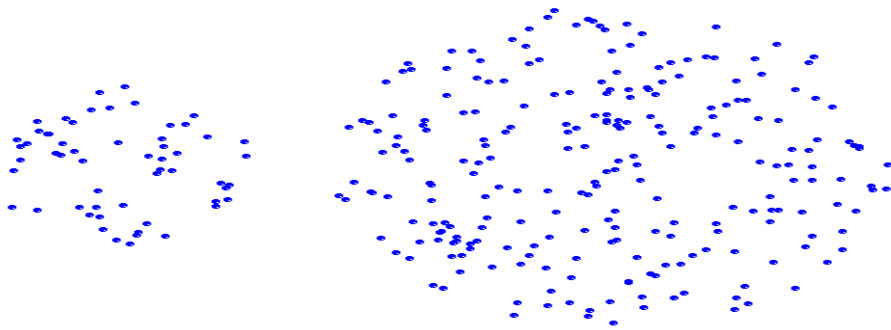


**Nested Clusters**

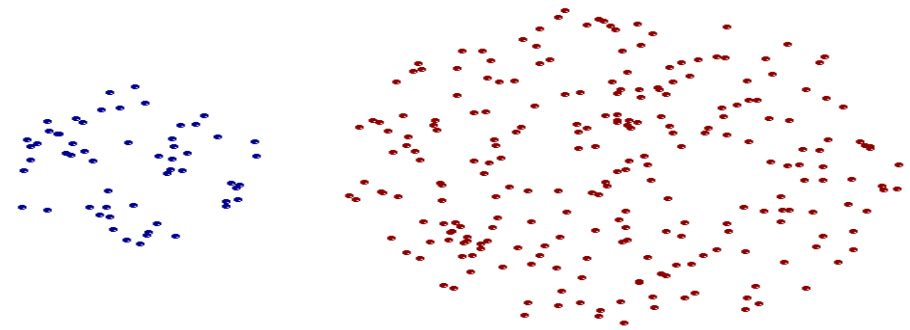


**Dendrogram**

# Strength of MIN



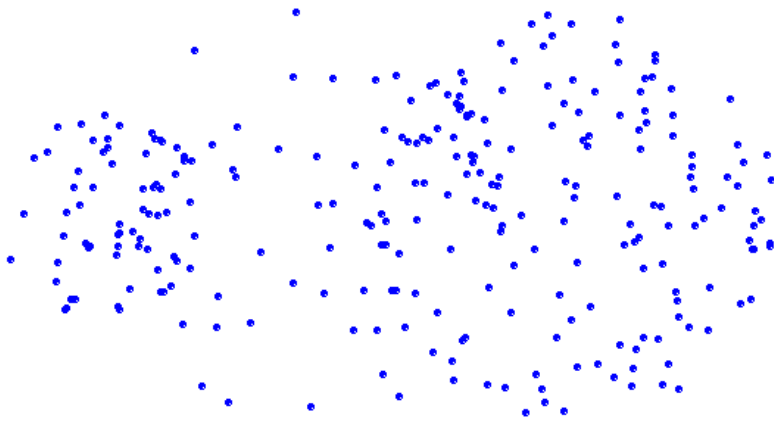
ORIGINAL POINTS



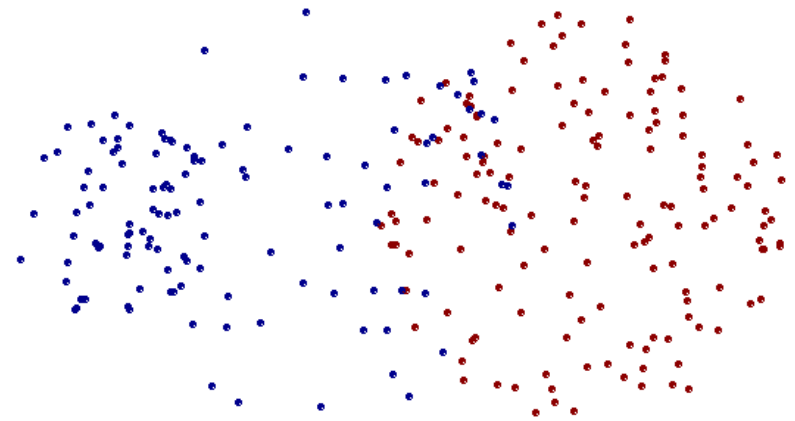
TWO CLUSTERS

CAN HANDLE NON-ELLIPTICAL SHAPES

# Limitations of MIN



ORIGINAL POINTS



TWO CLUSTERS

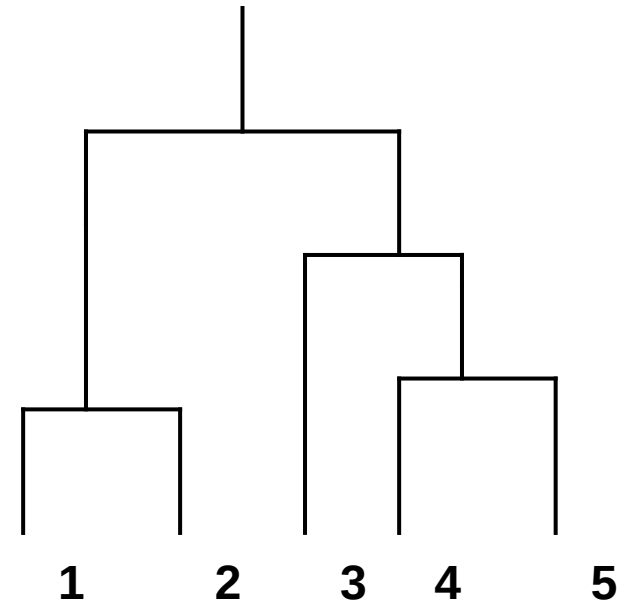
SENSITIVE TO NOISE AND OUTLIERS



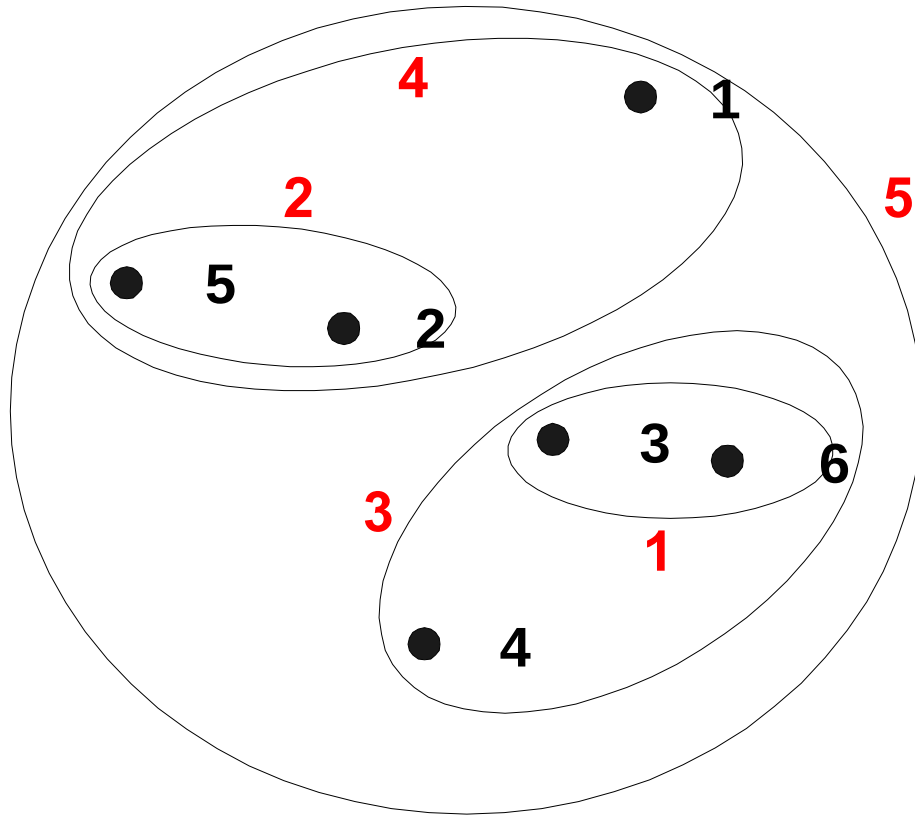
# Cluster Similarity: MAX or Complete Linkage

- ✗ SIMILARITY OF TWO CLUSTERS IS BASED ON THE TWO LEAST SIMILAR (MOST DISTANT) POINTS IN THE DIFFERENT CLUSTERS
- ✗ Determined by all pairs of points in the two clusters

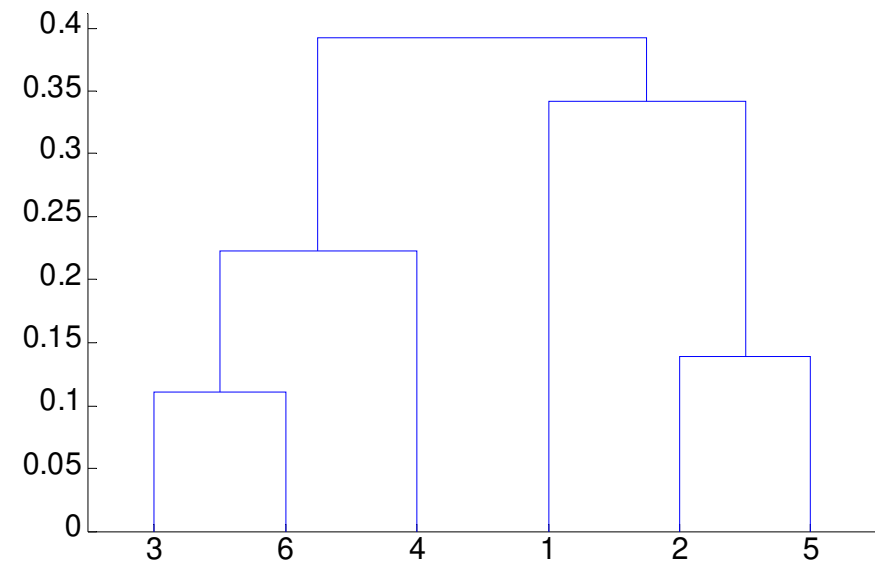
	I1	I2	I3	I4	I5
I1	1.00	0.90	0.10	0.65	0.20
I2	0.90	1.00	0.70	0.60	0.50
I3	0.10	0.70	1.00	0.40	0.30
I4	0.65	0.60	0.40	1.00	0.80
I5	0.20	0.50	0.30	0.80	1.00



# Hierarchical Clustering: MAX

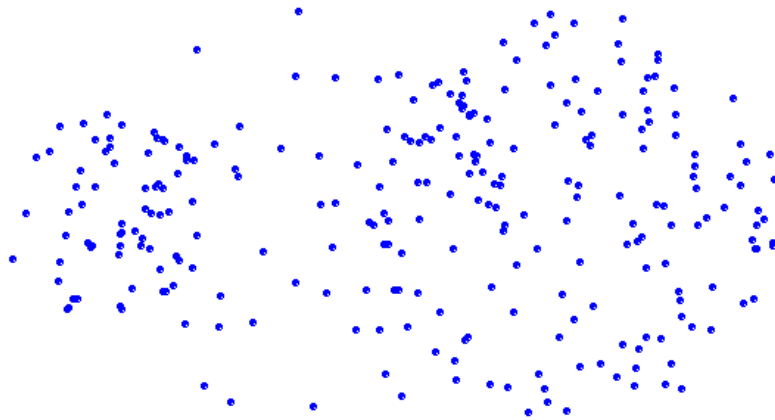


NESTED CLUSTERS

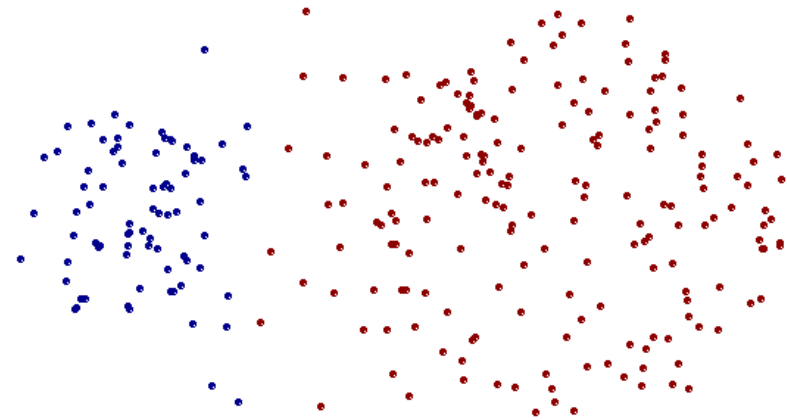


DENDROGRAM

# Strength of MAX



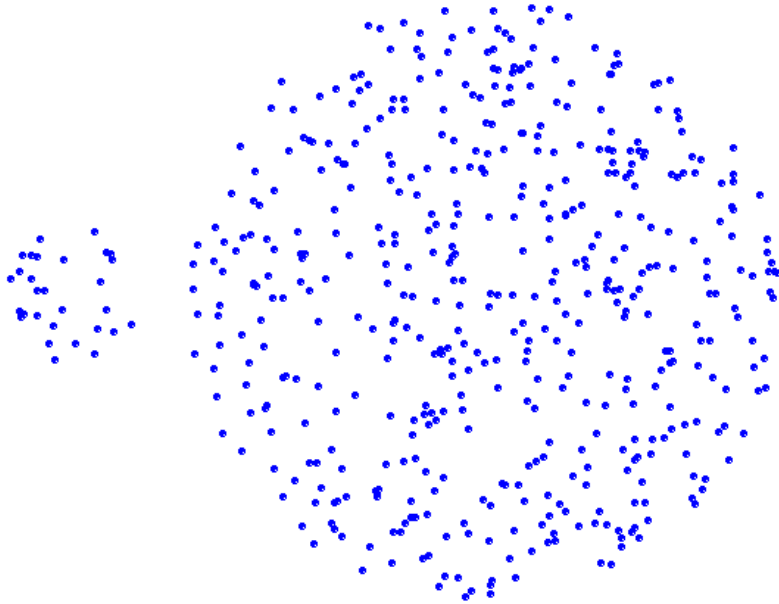
ORIGINAL POINTS



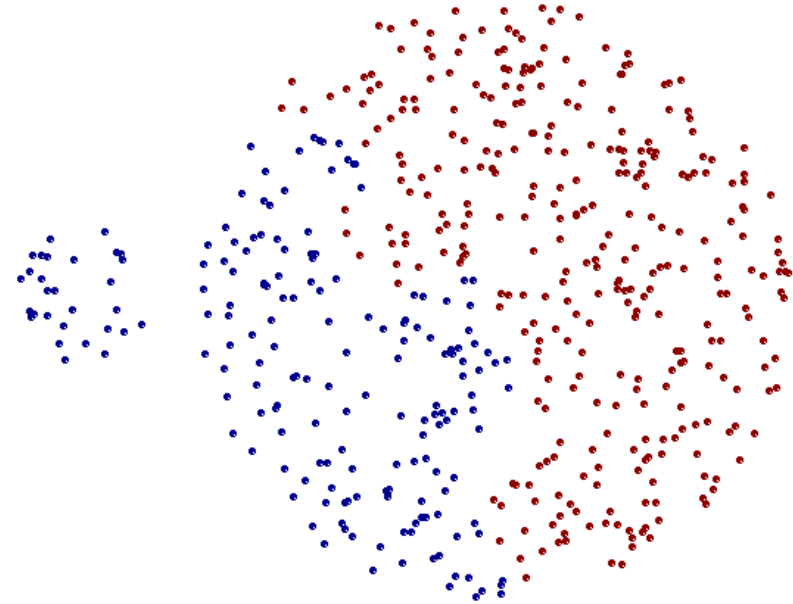
TWO CLUSTERS

LESS SUSCEPTIBLE TO NOISE AND OUTLIERS

# Limitations of MAX



ORIGINAL POINTS



TWO CLUSTERS

TENDS TO BREAK LARGE CLUSTERS

BIASED TOWARDS GLOBULAR CLUSTERS

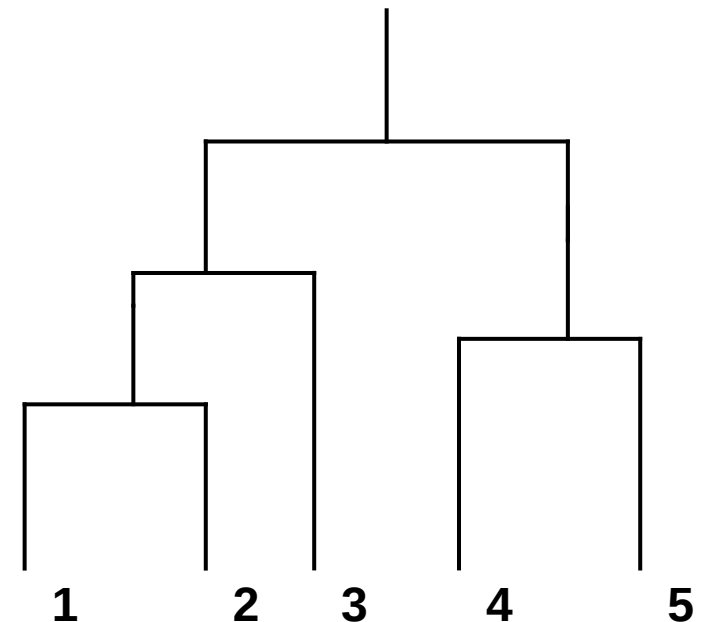
# Cluster Similarity: Group Average

✕ PROXIMITY OF TWO CLUSTERS IS THE AVERAGE OF PAIRWISE PROXIMITY BETWEEN POINTS IN THE TWO CLUSTERS.

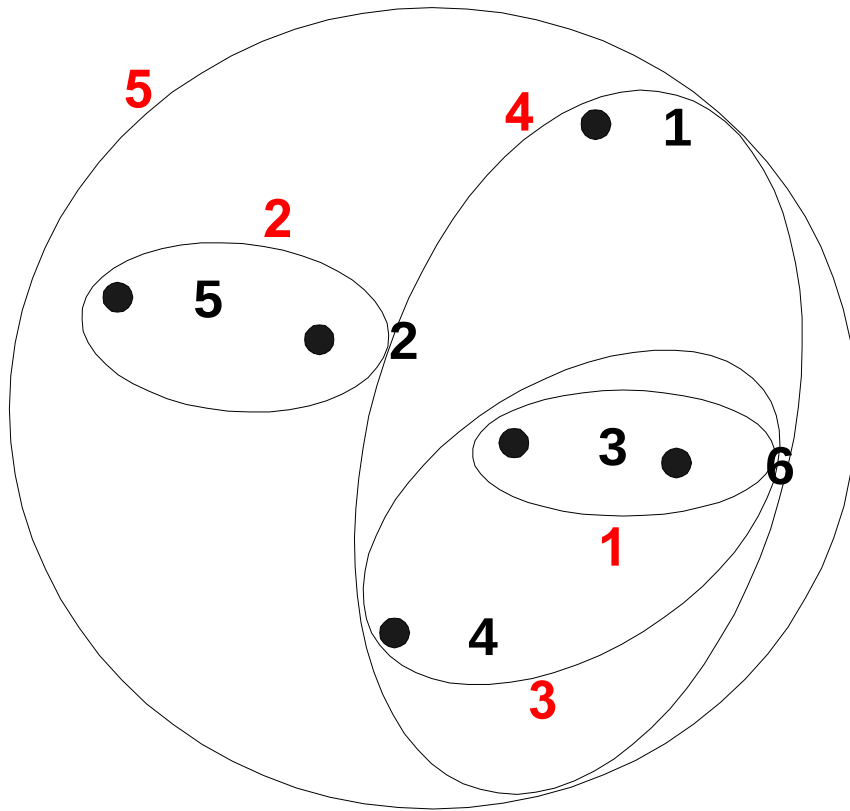
$$proximity(C_i, C_j) = \frac{\sum_{p_i \in C_i, p_j \in C_j} proximity(p_i, p_j)}{|C_i| \cdot |C_j|}$$

✕ NEED TO USE AVERAGE CONNECTIVITY FOR SCALABILITY SINCE TOTAL PROXIMITY FAVORS LARGE CLUSTERS

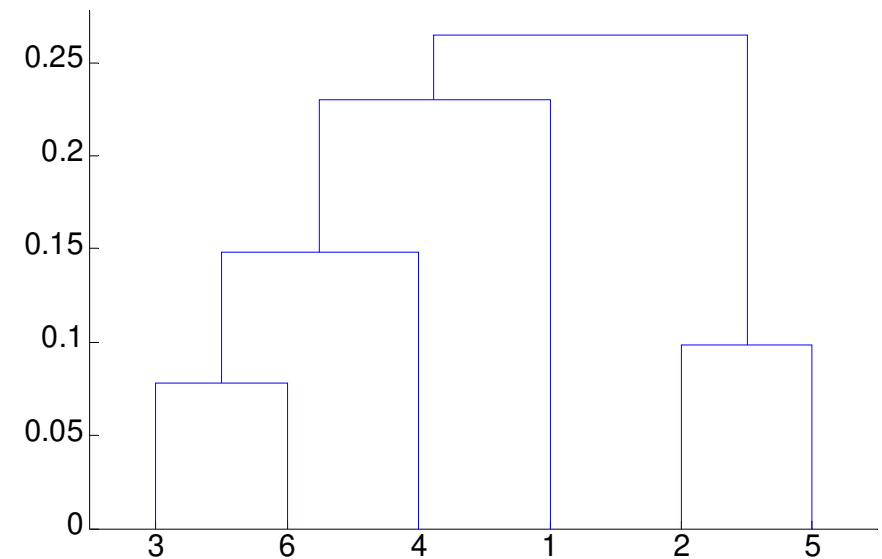
	I1	I2	I3	I4	I5
I1	1.00	0.90	0.10	0.65	0.20
I2	0.90	1.00	0.70	0.60	0.50
I3	0.10	0.70	1.00	0.40	0.30
I4	0.65	0.60	0.40	1.00	0.80
I5	0.20	0.50	0.30	0.80	1.00



# Hierarchical Clustering: Group Average



NESTED CLUSTERS



DENDROGRAM

# Hierarchical Clustering: Group Average

- ✗ COMPROMISE BETWEEN SINGLE AND COMPLETE LINK

- ✗ STRENGTHS

  - ✗ Less susceptible to noise and outliers

- ✗ LIMITATIONS

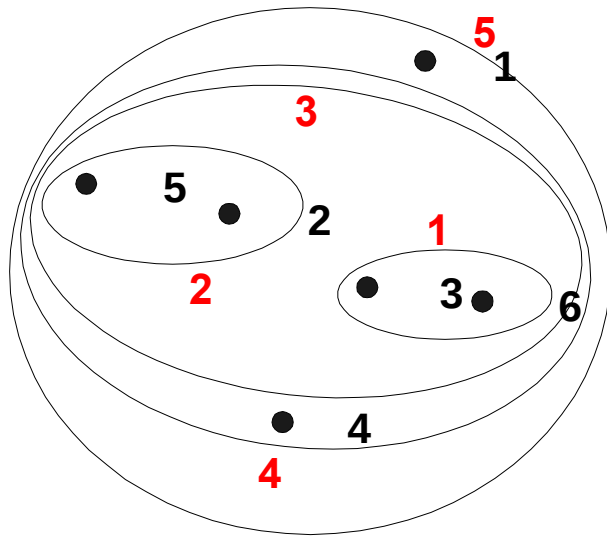
  - ✗ Biased towards globular clusters

# Cluster Similarity: Ward's Method

- ✗ SIMILARITY OF TWO CLUSTERS IS BASED ON THE INCREASE IN SQUARED ERROR WHEN TWO CLUSTERS ARE MERGED
  - ✗ Similar to group average if distance between points is distance squared
- ✗ LESS SUSCEPTIBLE TO NOISE AND OUTLIERS
- ✗ BIASED TOWARDS GLOBULAR CLUSTERS
- ✗ HIERARCHICAL ANALOGUE OF K-MEANS
  - ✗ Can be used to initialize K-means

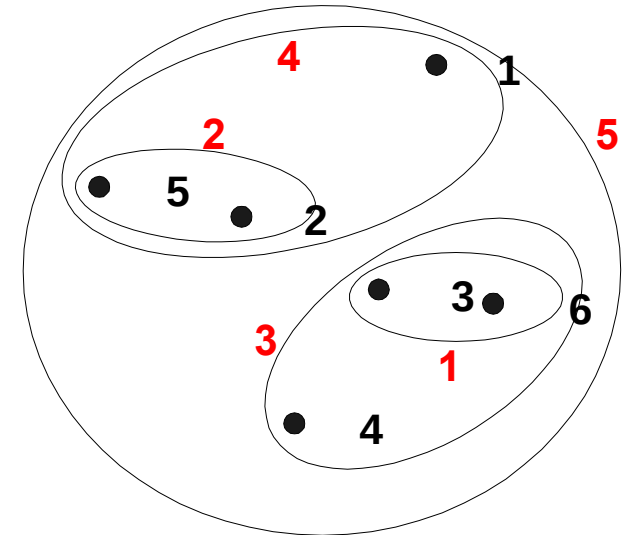


# Hierarchical Clustering: Comparison



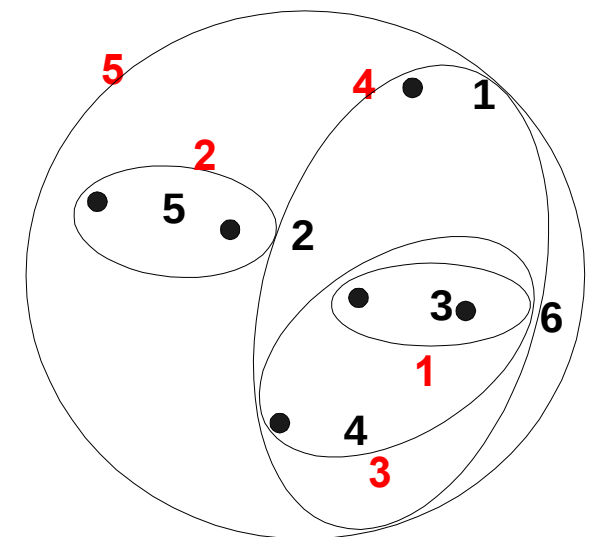
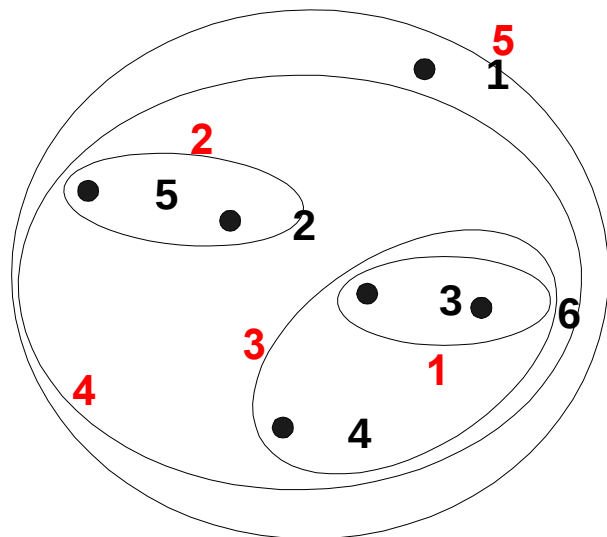
MIN

MAX



WARD'S METHOD

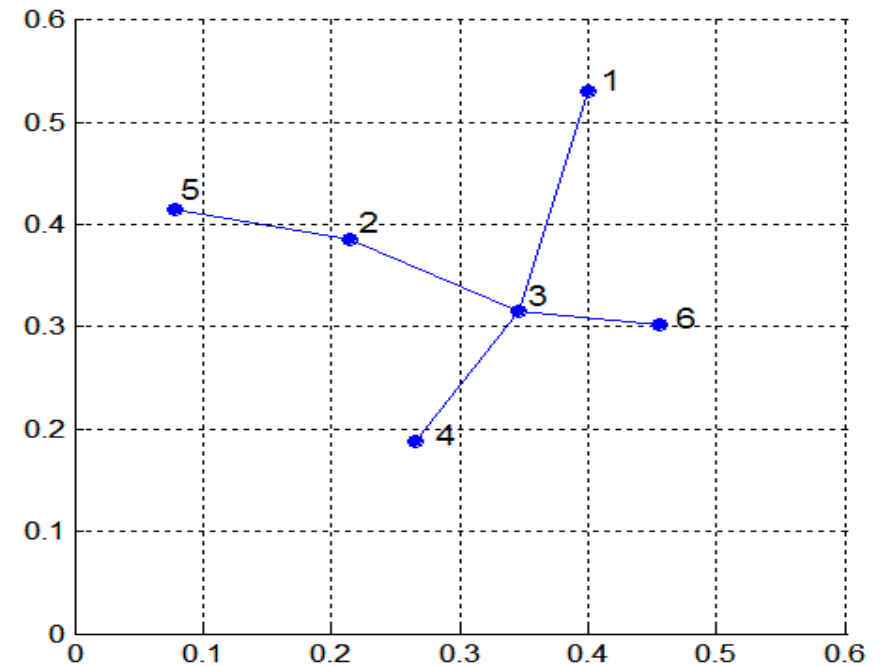
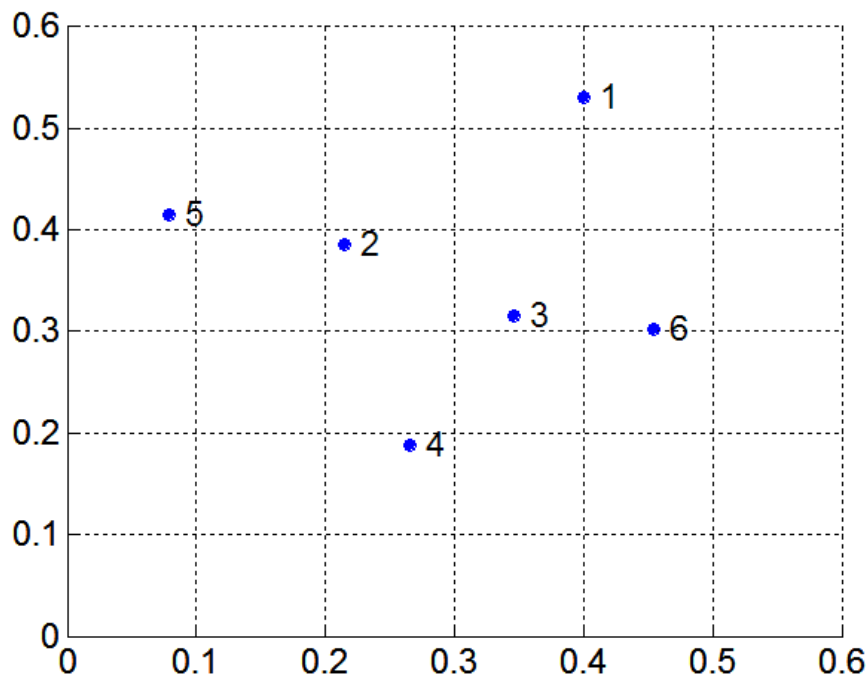
GROUP AVERAGE



# MST: Divisive Hierarchical Clustering

## ✕ BUILD MST (MINIMUM SPANNING TREE)

- ✕ Start with a tree that consists of any point
- ✕ In successive steps, look for the closest pair of points (p, q) such that one point (p) is in the current tree but the other (q) is not
- ✕ Add q to the tree and put an edge between p and q



# MST: Divisive Hierarchical Clustering

✕USE MST FOR CONSTRUCTING HIERARCHY OF CLUSTERS

---

**Algorithm 7.5** MST Divisive Hierarchical Clustering Algorithm

---

- 1: Compute a minimum spanning tree for the proximity graph.
  - 2: **repeat**
  - 3:   Create a new cluster by breaking the link corresponding to the largest distance (smallest similarity).
  - 4: **until** Only singleton clusters remain
-

# Hierarchical Clustering: Time and Space requirements

- ✗  $O(N^2)$  SPACE SINCE IT USES THE PROXIMITY MATRIX.

  - ✗  $N$  is the number of points.

- ✗  $O(N^3)$  TIME IN MANY CASES

  - ✗ There are  $N$  steps and at each step the size,  $N^2$ , proximity matrix 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

# Example of Hierarchical Clustering using String Edit Distance

## Pedro (Portuguese)

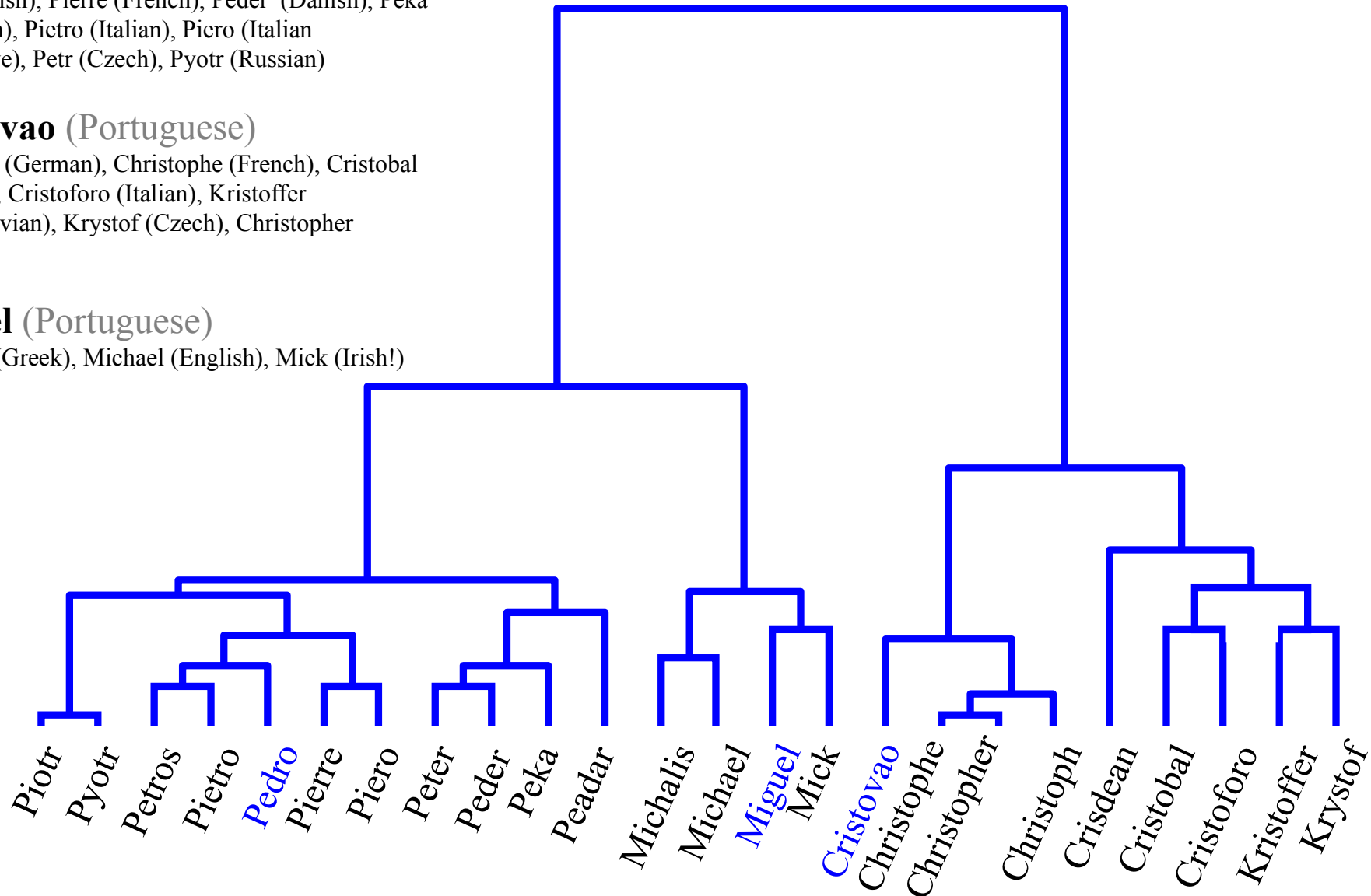
Petros (Greek), Peter (English), Piotr (Polish), Peadar (Irish), Pierre (French), Peder (Danish), Peka (Hawaiian), Pietro (Italian), Piero (Italian Alternative), Petr (Czech), Pyotr (Russian)

## Cristovao (Portuguese)

Christoph (German), Christophe (French), Cristobal (Spanish), Cristoforo (Italian), Kristoffer (Scandinavian), Krystof (Czech), Christopher (English)

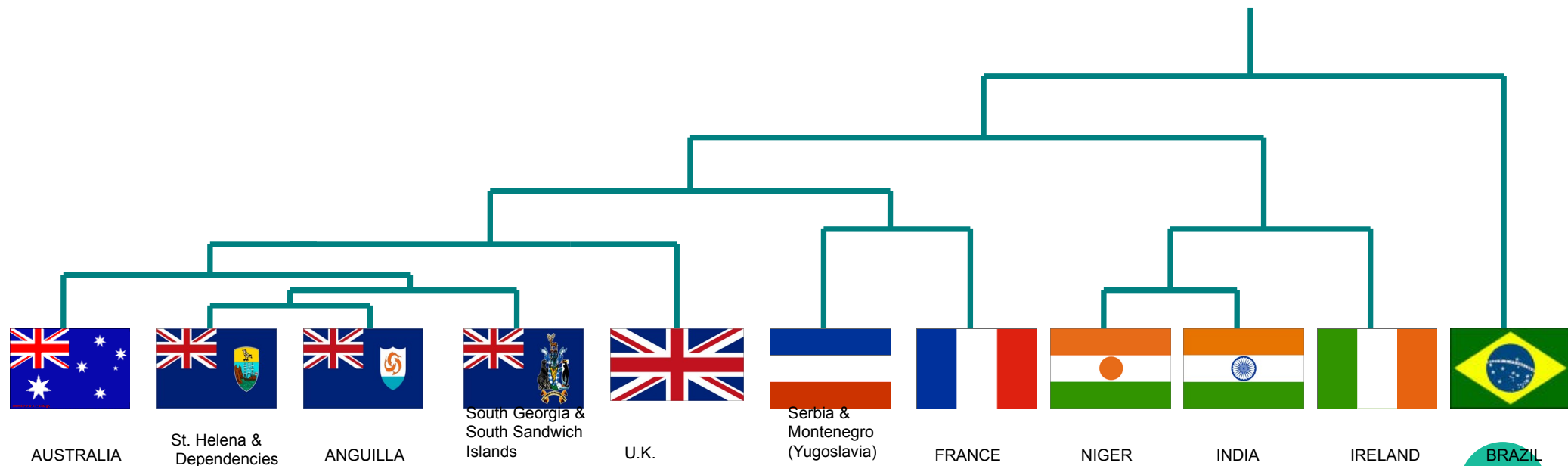
## Miguel (Portuguese)

Michalis (Greek), Michael (English), Mick (Irish!)

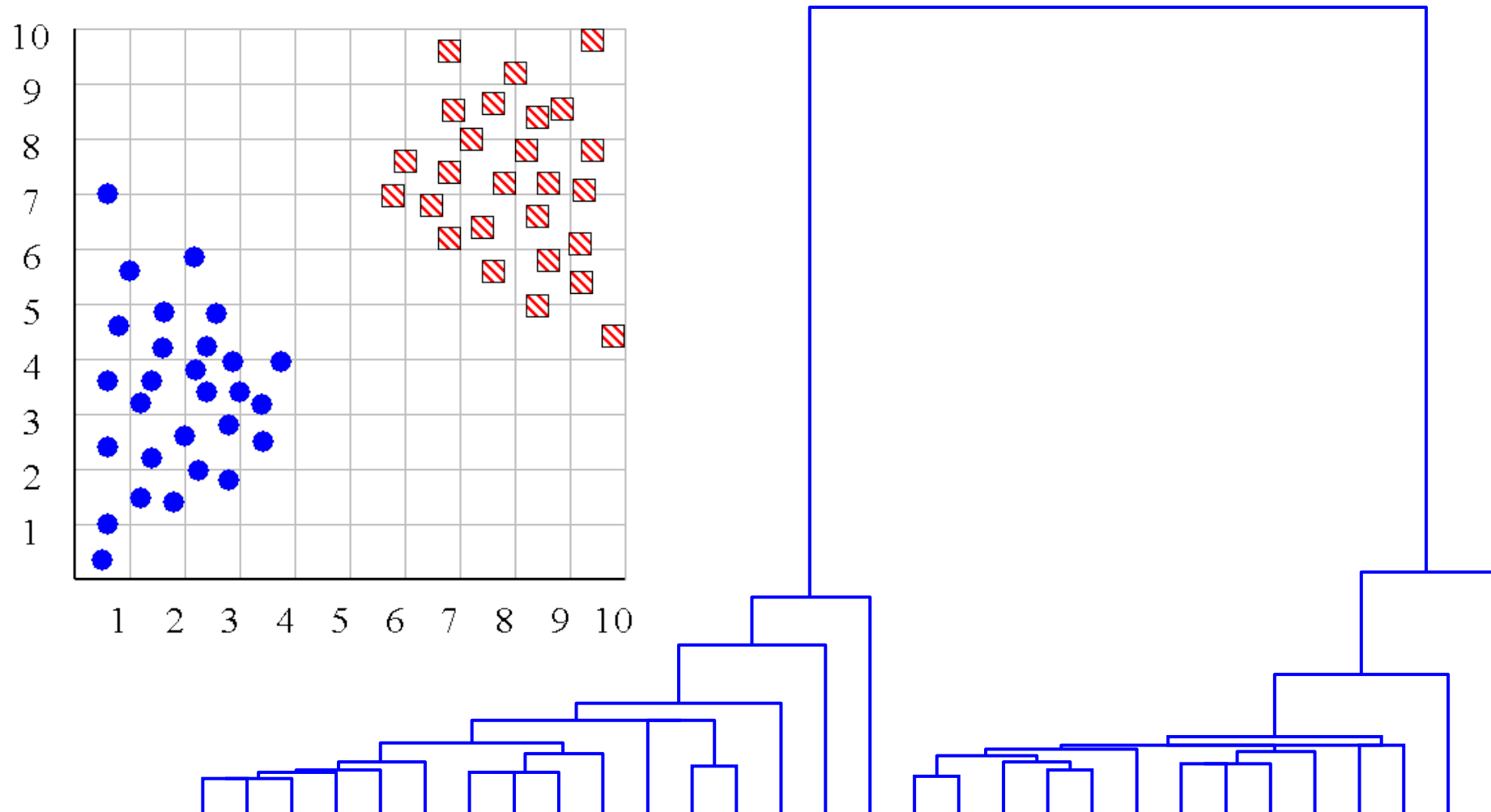


# Spurious clusters

- ✗ HIERARCHICAL CLUSTERING CAN SOMETIMES SHOW PATTERNS THAT ARE MEANINGLESS OR SPURIOUS
- ✗ THE TIGHT GROUPING OF AUSTRALIA, ANGUILLA, ST. HELENA ETC IS MEANINGFUL; ALL THESE COUNTRIES ARE FORMER UK COLONIES
- ✗ HOWEVER THE TIGHT GROUPING OF NIGER AND INDIA IS COMPLETELY SPURIOUS; THERE IS NO CONNECTION BETWEEN THE TWO.



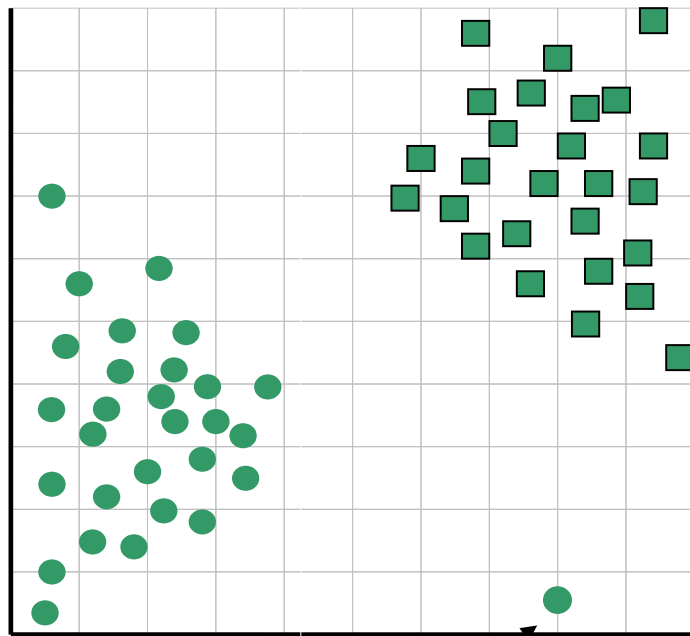
We can use the dendrogram to determine the “correct” number of clusters.



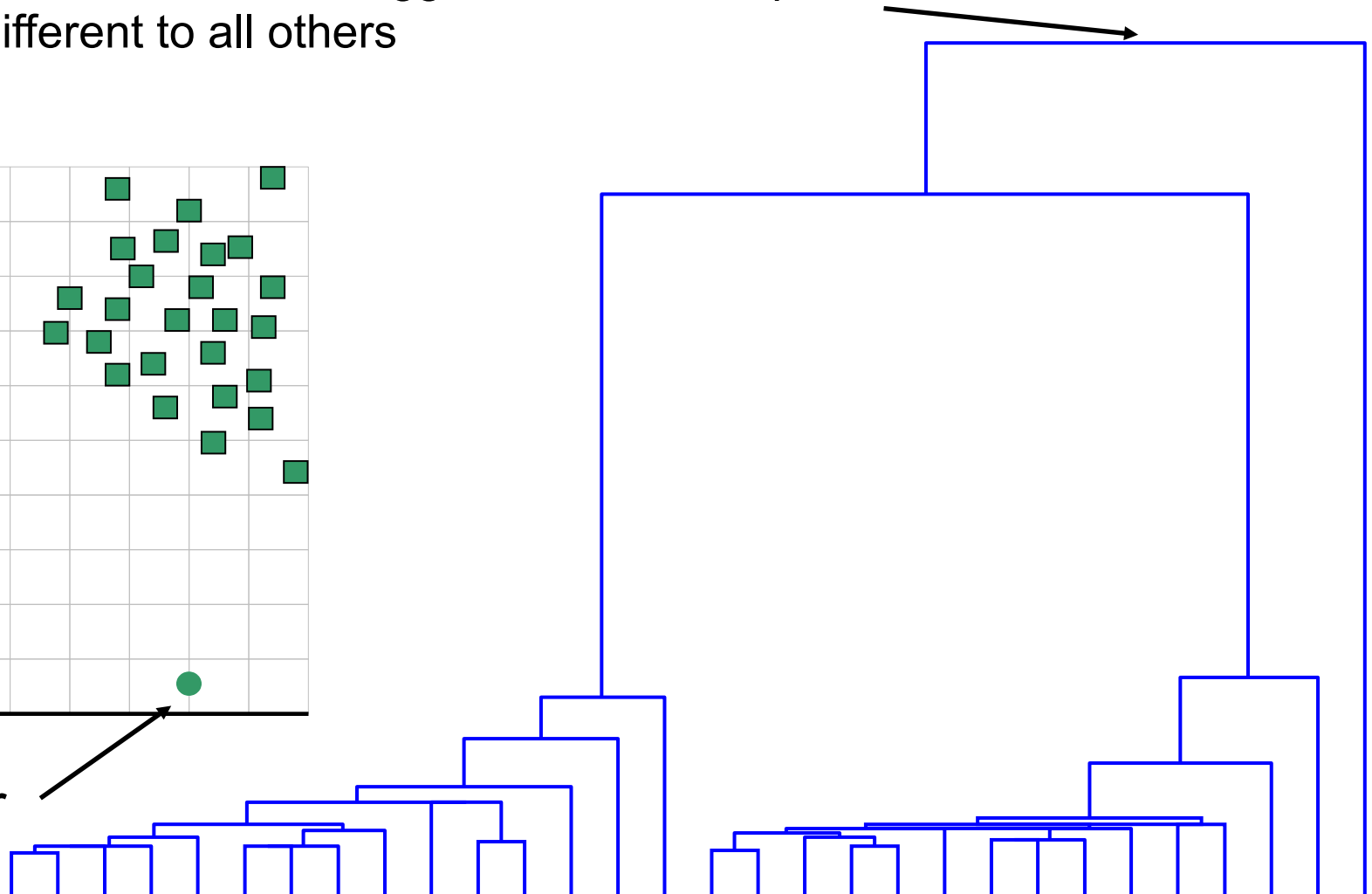


# One potential use of a dendrogram: detecting outliers

The single isolated branch is suggestive of a data point that is very different to all others

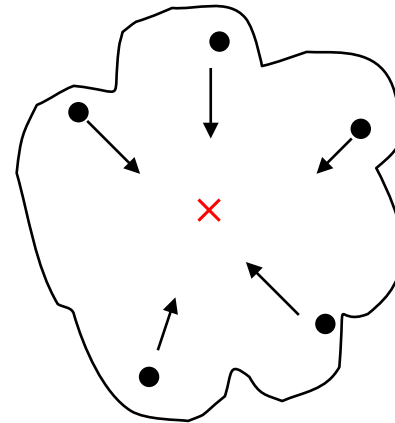
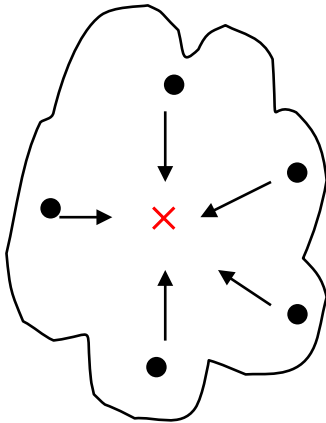


Outlier



# CURE: Another Hierarchical Approach

✗ USES A NUMBER OF POINTS TO REPRESENT A CLUSTER



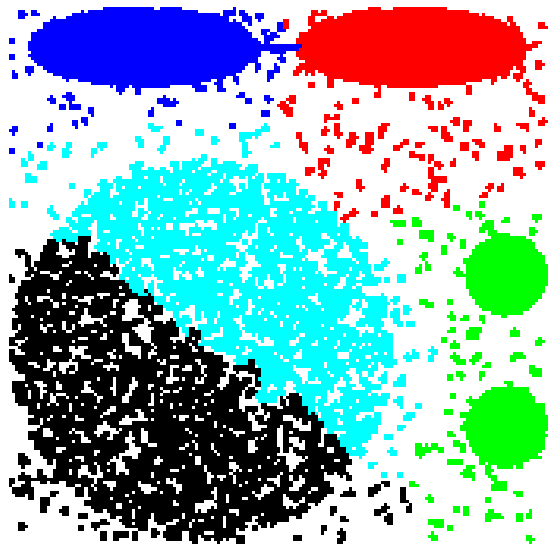
✗ REPRESENTATIVE POINTS ARE FOUND BY SELECTING A CONSTANT NUMBER OF POINTS FROM A CLUSTER AND THEN “SHRINKING” THEM TOWARD THE CENTER OF THE CLUSTER

✗ CLUSTER SIMILARITY IS THE SIMILARITY OF THE CLOSEST PAIR OF REPRESENTATIVE POINTS FROM DIFFERENT CLUSTERS

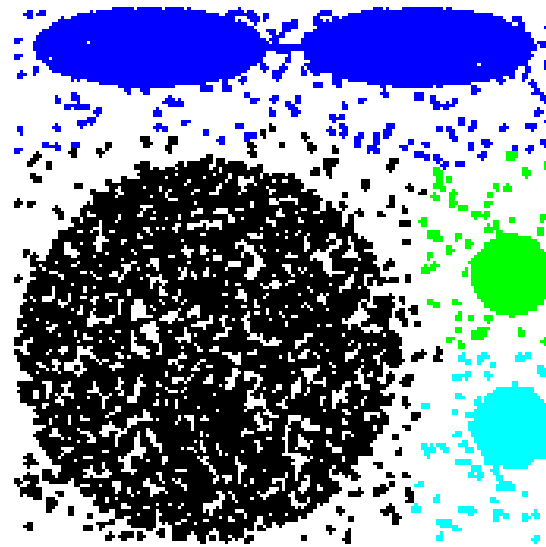
# CURE

- ✗ SHRINKING REPRESENTATIVE POINTS TOWARD THE CENTER HELPS AVOID PROBLEMS WITH NOISE AND OUTLIERS
- ✗ CURE IS BETTER ABLE TO HANDLE CLUSTERS OF ARBITRARY SHAPES AND SIZES

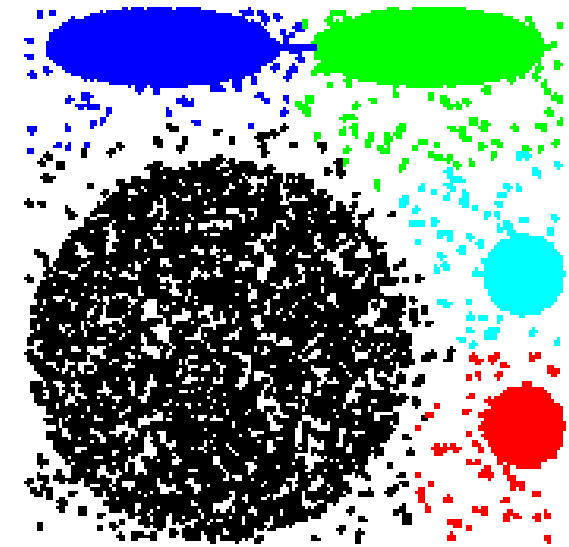
# Experimental Results: CURE



a) BIRCH

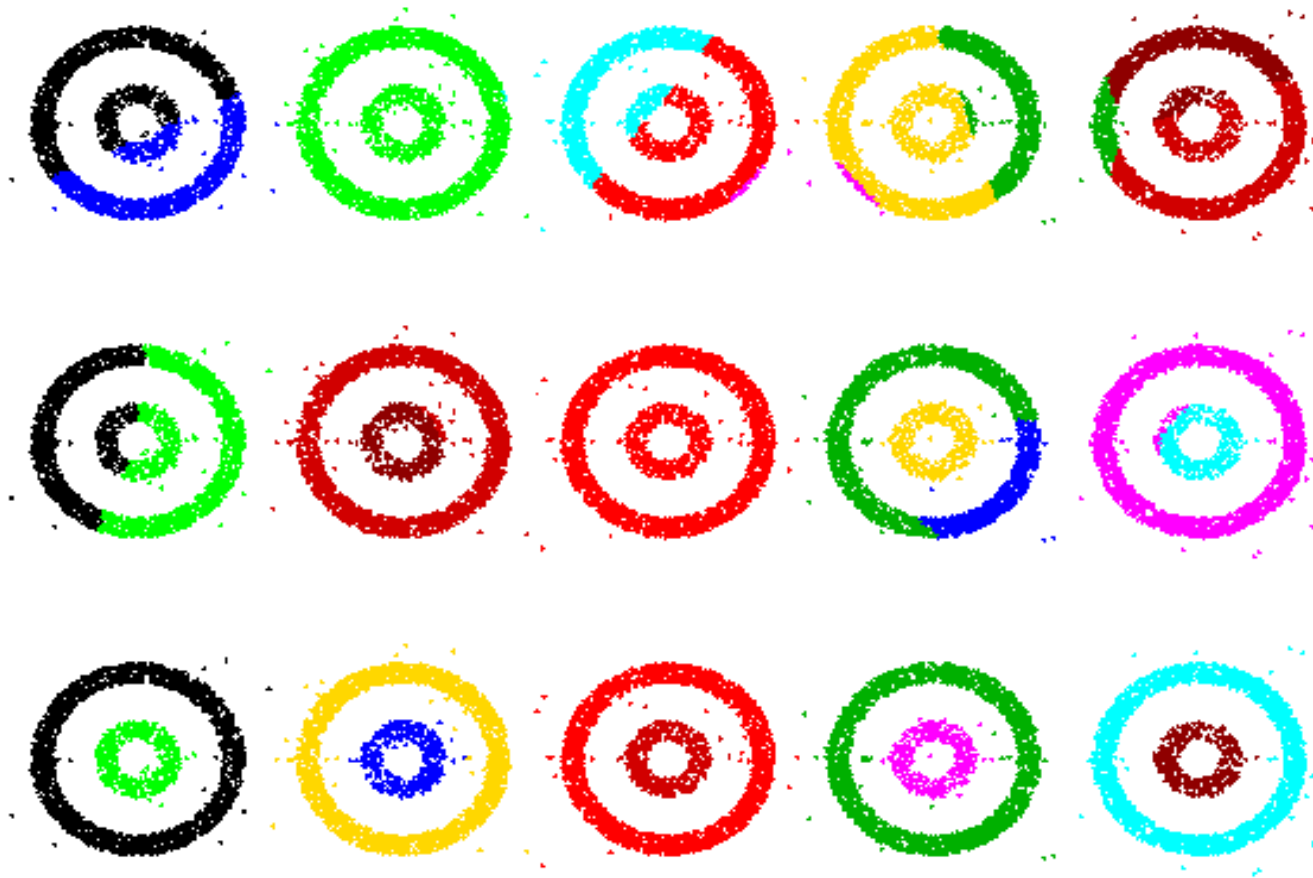


b) MST METHOD



c) CURE

# Experimental Results: CURE



a) BIRCH

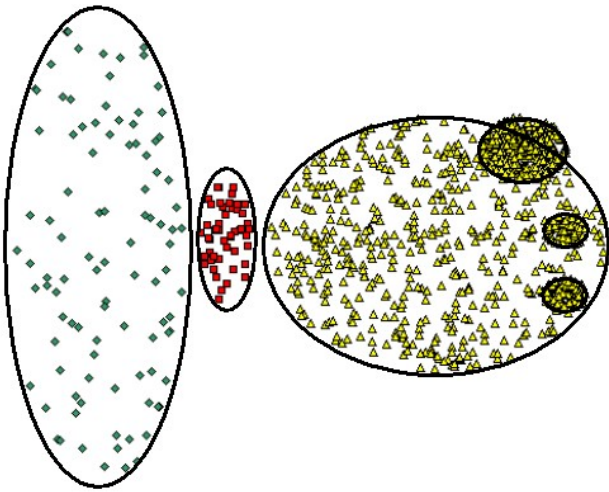
(centroid)

b) MST METHOD

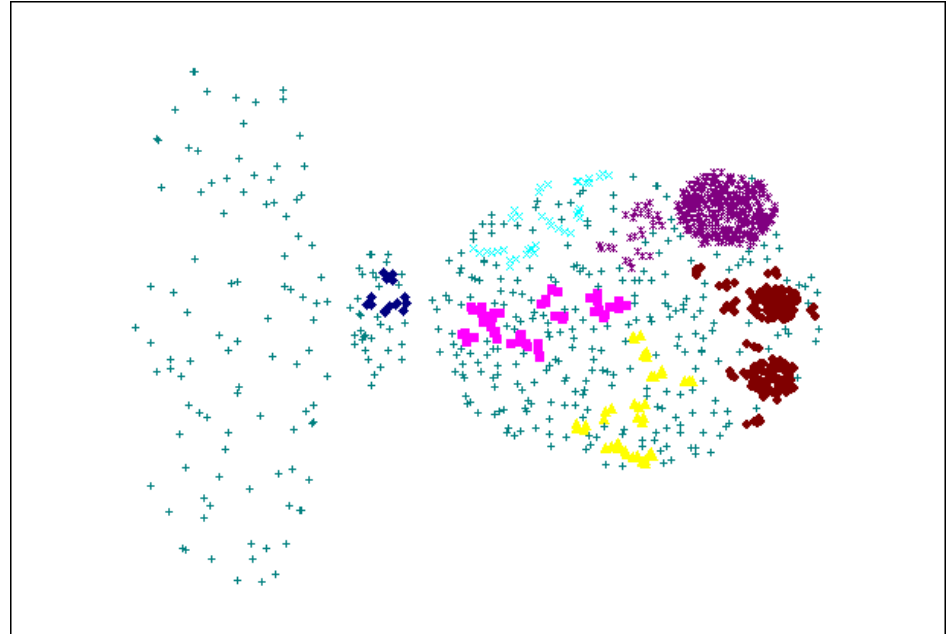
(single link)

c) CURE

# CURE Cannot Handle Differing Densities



**Original Points**



**CURE**

# Recent Hierarchical Clustering Methods

## ✗ MAJOR WEAKNESS OF AGGLOMERATIVE CLUSTERING METHODS

- ✗ do not scale well: time complexity of at least  $O(n^2)$ , where  $n$  is the number of total objects
- ✗ can never undo what was done previously

## ✗ INTEGRATION OF HIERARCHICAL WITH DISTANCE-BASED CLUSTERING

- ✗ BIRCH: uses CF-tree and incrementally adjusts the quality of sub-clusters
- ✗ ROCK: clustering categorical data by neighbor and link analysis
- ✗ CHAMELEON: hierarchical clustering using dynamic modeling

# BIRCH (1996)

- ✗ BIRCH: BALANCED ITERATIVE REDUCING AND CLUSTERING USING HIERARCHIES (ZHANG, RAMAKRISHNAN & LIVNY, SIGMOD'96)
- ✗ INCREMENTALLY CONSTRUCT A CF (CLUSTERING FEATURE) TREE, A HIERARCHICAL DATA STRUCTURE FOR MULTIPHASE CLUSTERING
  - ✗ Phase 1: scan DB to build an initial in-memory CF tree (a multi-level compression of the data that tries to preserve the inherent clustering structure of the data)
  - ✗ Phase 2: use an arbitrary clustering algorithm to cluster the leaf nodes of the CF-tree
- ✗ *SCALES LINEARLY*: FINDS A GOOD CLUSTERING WITH A SINGLE SCAN AND IMPROVES THE QUALITY WITH A FEW ADDITIONAL SCANS
- ✗ *WEAKNESS*: HANDLES ONLY NUMERIC DATA, AND SENSITIVE TO THE ORDER OF THE DATA RECORD.



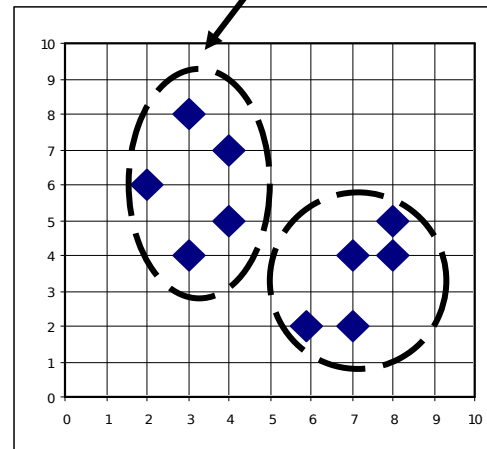
# Clustering Feature Vector in BIRCH

**Clustering Feature:**  $CF = (N, LS, \overrightarrow{SS})$

$N$ : **Number of data points**

$LS$ :  $\sum_{i=1}^N X_i \rightarrow$

$SS$ :  $\sum_{i=1}^N X_i^2 \rightarrow$



$CF = (5, 16.30, 54.190)$

(3,4)

(2,6)

(4,5)

(4,7)

(3,8)

# CF-Tree in BIRCH

## ✕CLUSTERING FEATURE:

- ✕summary of the statistics for a given subcluster: the 0-th, 1st and 2nd moments of the subcluster from the statistical point of view.
- ✕registers crucial measurements for computing cluster and utilizes storage efficiently

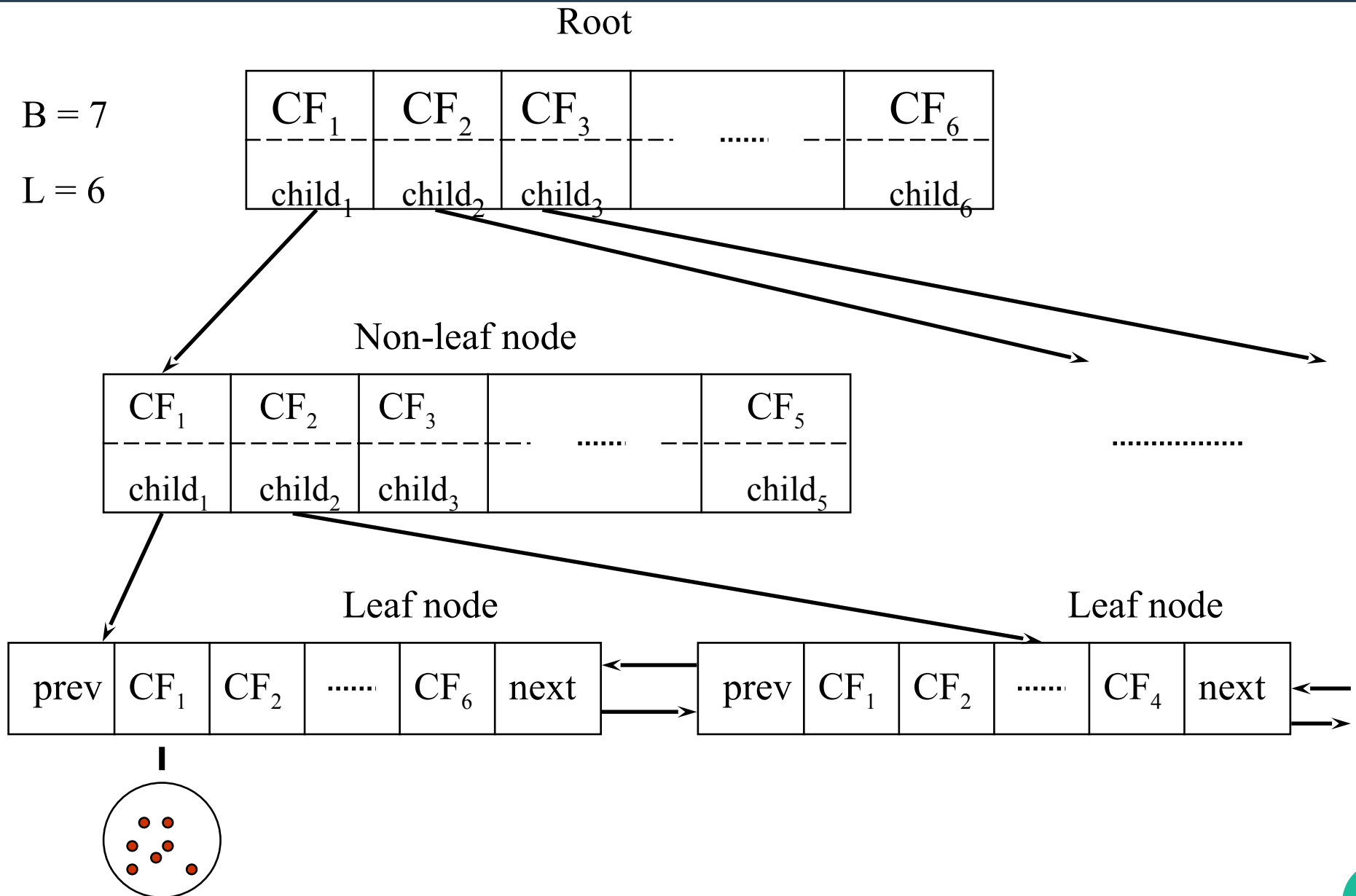
## ■ A CF TREE IS A HEIGHT-BALANCED TREE THAT STORES THE CLUSTERING FEATURES FOR A HIERARCHICAL CLUSTERING

- ✕A nonleaf node in a tree has descendants or “children”
- ✕The nonleaf nodes store sums of the CFs of their children

## ✕A CF TREE HAS TWO PARAMETERS

- ✕Branching factor: specify the maximum number of children
- ✕threshold: max diameter of sub-clusters stored at the leaf nodes

# The CF Tree Structure



# Clustering Categorical Data: The ROCK Algorithm

- ✕ ROCK: ROBUST CLUSTERING USING LINKS

  - ✕ S. Guha, R. Rastogi & K. Shim, ICDE'99

- ✕ MAJOR IDEAS

  - ✕ Use links to measure similarity/proximity

  - ✕ Not distance-based

  - ✕ Computational complexity:  $O(n^2 + nm_m m_a + n^2 \log n)$

- ✕ ALGORITHM: SAMPLING-BASED CLUSTERING

  - ✕ Draw random sample

  - ✕ Cluster with links

  - ✕ Label data in disk

- ✕ EXPERIMENTS

  - ✕ Congressional voting, mushroom data

# Similarity Measure in ROCK

- TRADITIONAL MEASURES FOR CATEGORICAL DATA MAY NOT WORK WELL, E.G., JACCARD COEFFICIENT
- EXAMPLE: TWO GROUPS (CLUSTERS) OF TRANSACTIONS
  - $C_1$ .  $\langle a, b, c, d, e \rangle$ :  $\{a, b, c\}, \{a, b, d\}, \{a, b, e\}, \{a, c, d\}, \{a, c, e\}, \{a, d, e\}, \{b, c, d\}, \{b, c, e\}, \{b, d, e\}, \{c, d, e\}$
  - $C_2$ .  $\langle a, b, f, g \rangle$ :  $\{a, b, f\}, \{a, b, g\}, \{a, f, g\}, \{b, f, g\}$
- JACCARD CO-EFFICIENT MAY LEAD TO WRONG CLUSTERING RESULT
  - $C_1$ : 0.2 ( $\{a, b, c\}, \{b, d, e\}$ ) to 0.5 ( $\{a, b, c\}, \{a, b, d\}$ )
  - $C_1$  &  $C_2$ : could be as high as 0.5 ( $\{a, b, c\}, \{a, b, f\}$ )
- JACCARD CO-EFFICIENT-BASED SIMILARITY FUNCTION:  $Sim(T_1, T_2) = \frac{|T_1 \cap T_2|}{|T_1 \cup T_2|}$
- Ex. Let  $T_1 = \{a, b, c\}, T_2 = \{c, d, e\}$

$$Sim(T_1, T_2) = \frac{|\{c\}|}{|\{a, b, c, d, e\}|} = \frac{1}{5} = 0.2$$

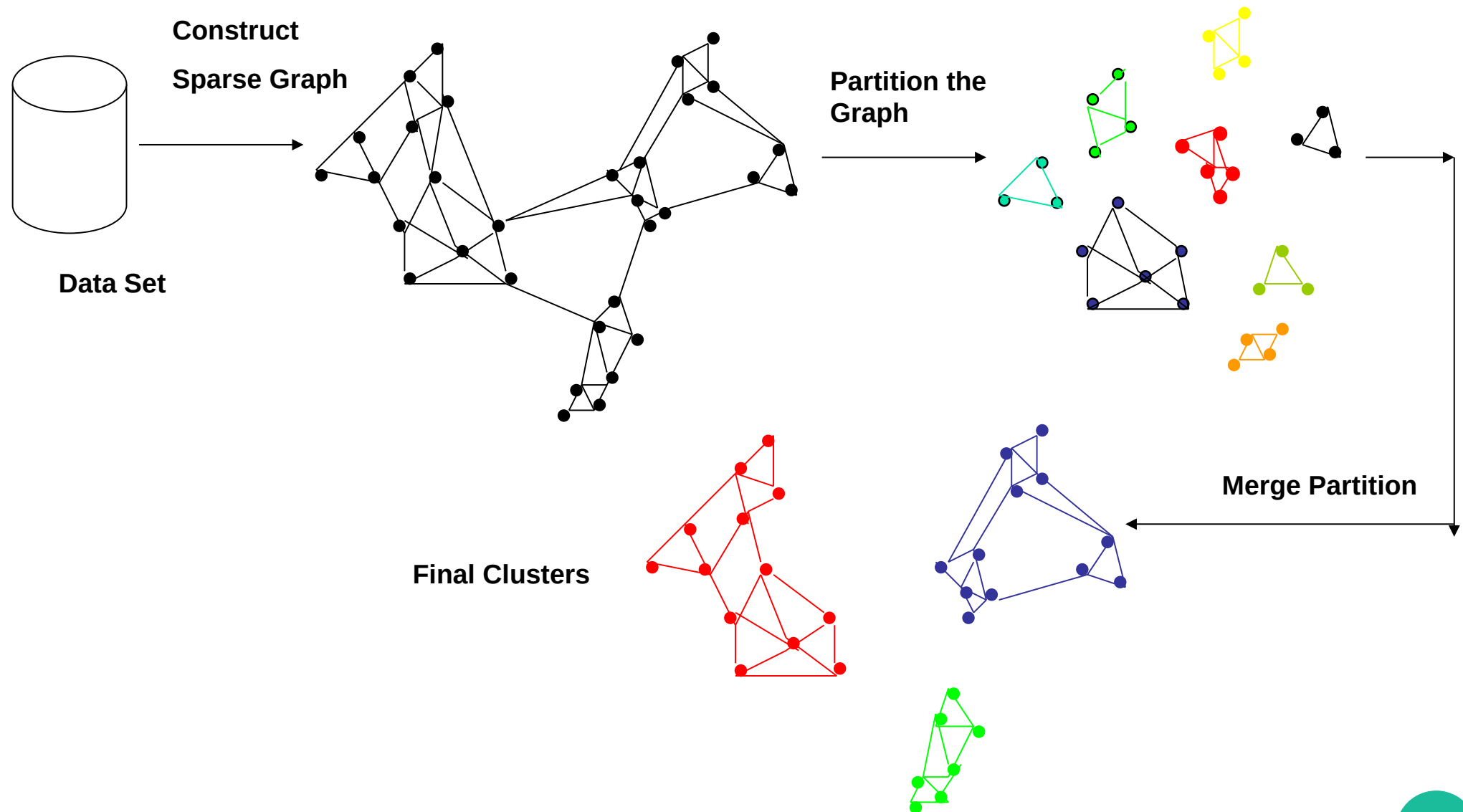
# Link Measure in ROCK

- LINKS: # OF COMMON NEIGHBORS
  - $C_1 \langle a, b, c, d, e \rangle$ :  $\{a, b, c\}$ ,  $\{a, b, d\}$ ,  $\{a, b, e\}$ ,  $\{a, c, d\}$ ,  $\{a, c, e\}$ ,  $\{a, d, e\}$ ,  $\{b, c, d\}$ ,  $\{b, c, e\}$ ,  $\{b, d, e\}$ ,  $\{c, d, e\}$
  - $C_2 \langle a, b, f, g \rangle$ :  $\{a, b, f\}$ ,  $\{a, b, g\}$ ,  $\{a, f, g\}$ ,  $\{b, f, g\}$
- LET  $T_1 = \{A, B, C\}$ ,  $T_2 = \{C, D, E\}$ ,  $T_3 = \{A, B, F\}$ 
  - $\text{link}(T_1, T_2) = 4$ , *since they have 4 common neighbors*
    - $\{a, c, d\}$ ,  $\{a, c, e\}$ ,  $\{b, c, d\}$ ,  $\{b, c, e\}$
  - $\text{link}(T_1, T_3) = 3$ , *since they have 3 common neighbors*
    - $\{a, b, d\}$ ,  $\{a, b, e\}$ ,  $\{a, b, g\}$
- THUS LINK IS A BETTER MEASURE THAN JACCARD COEFFICIENT

# CHAMELEON: Hierarchical Clustering Using Dynamic Modeling

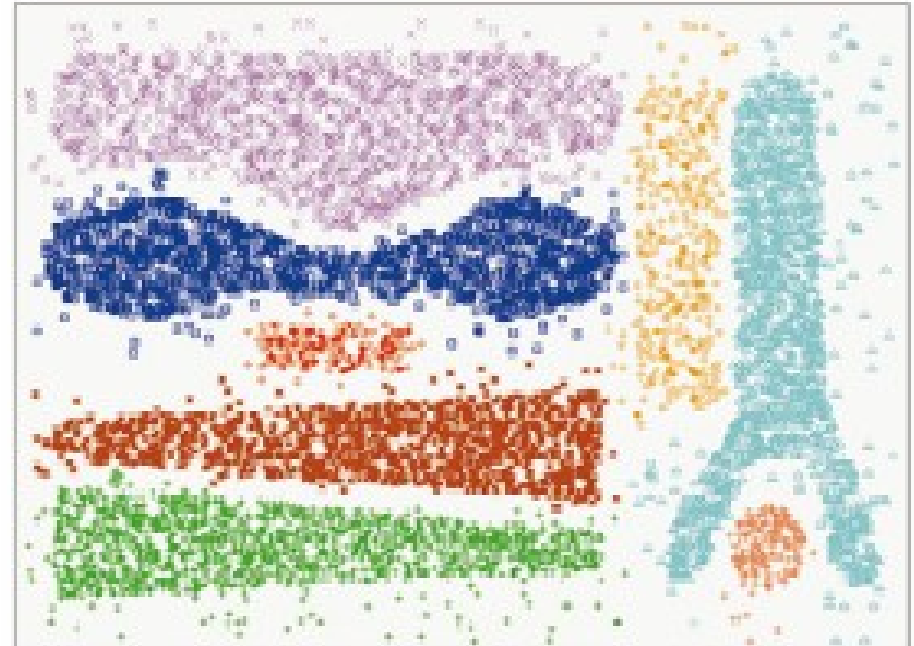
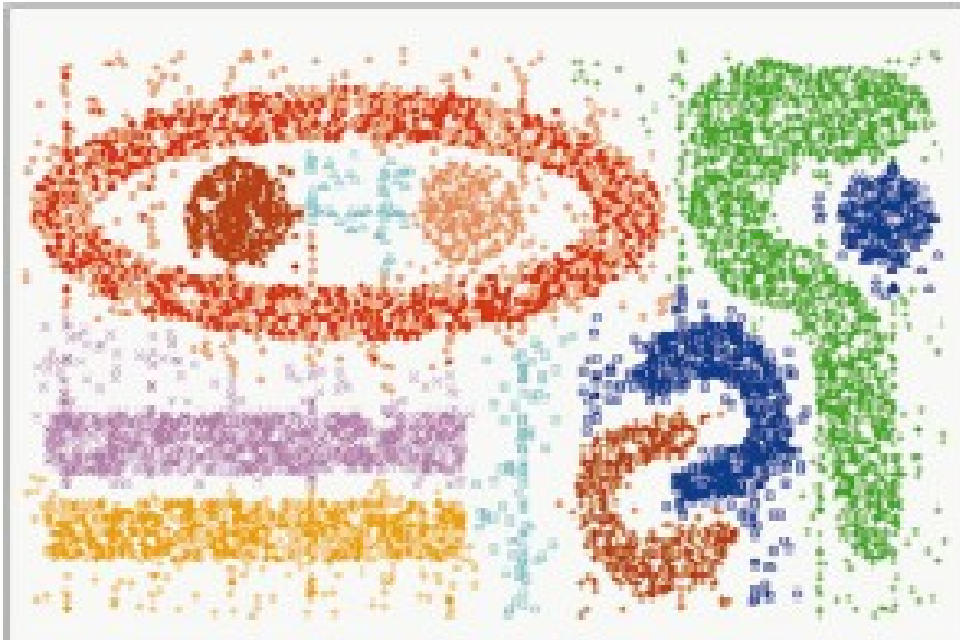
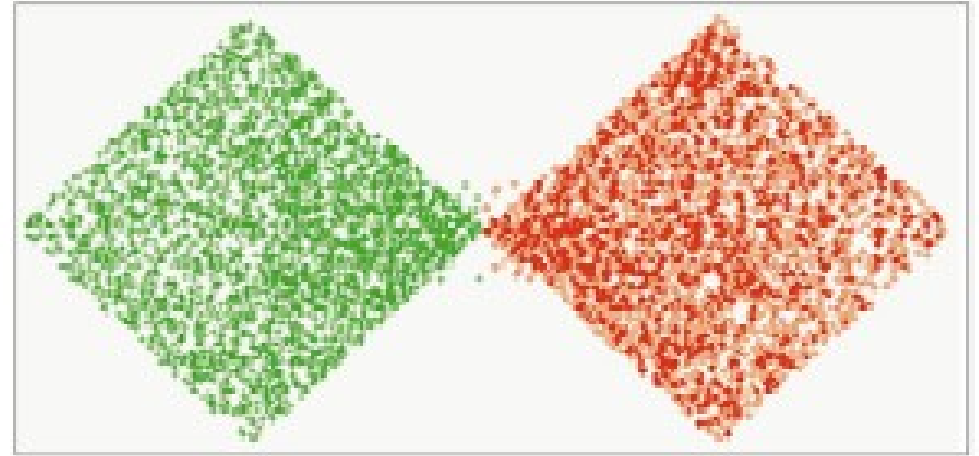
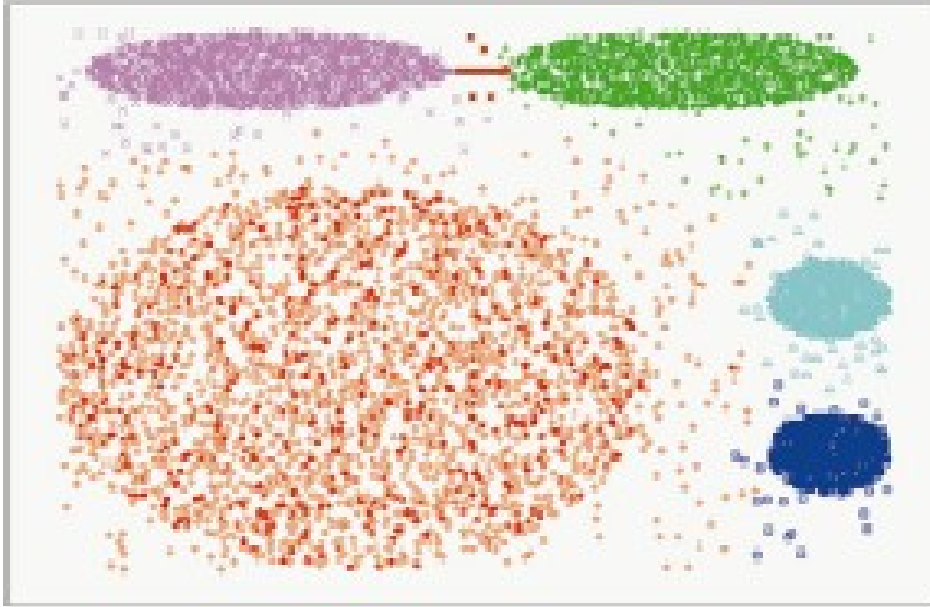
- CHAMELEON: BY G. KARYPIS, E.H. HAN, AND V. KUMAR'99
- MEASURES THE SIMILARITY BASED ON A DYNAMIC MODEL
  - Two clusters are merged only if the *interconnectivity* and *closeness (proximity)* between two clusters are high *relative to* the internal interconnectivity of the clusters and closeness of items within the clusters
  - **Cure** ignores information about **interconnectivity** of the objects, **Rock** ignores information about the **closeness** of two clusters
- A TWO-PHASE ALGORITHM
  1. Use a graph partitioning algorithm: cluster objects into a large number of relatively small sub-clusters
  2. Use an agglomerative hierarchical clustering algorithm: find the genuine clusters by repeatedly combining these sub-clusters

# Overall Framework of CHAMELEON





# CHAMELEON (Clustering Complex Objects)



# Graph-Based Clustering

- ✕ GRAPH-BASED CLUSTERING USES THE PROXIMITY GRAPH
  - ✕ Start with the proximity matrix
  - ✕ Consider each point as a node in a graph
  - ✕ Each edge between two nodes has a weight which is the proximity between the two points
  - ✕ Initially the proximity graph is fully connected
  - ✕ MIN (single-link) and MAX (complete-link) can be viewed as starting with this graph
- ✕ IN THE SIMPLEST CASE, CLUSTERS ARE CONNECTED COMPONENTS IN THE GRAPH

# Graph-Based Clustering: Sparsification

- THE AMOUNT OF DATA THAT NEEDS TO BE PROCESSED IS DRASTICALLY REDUCED
  - Sparsification can eliminate more than 99% of the entries in a proximity matrix
  - The amount of time required to cluster the data is drastically reduced
  - The size of the problems that can be handled is increased

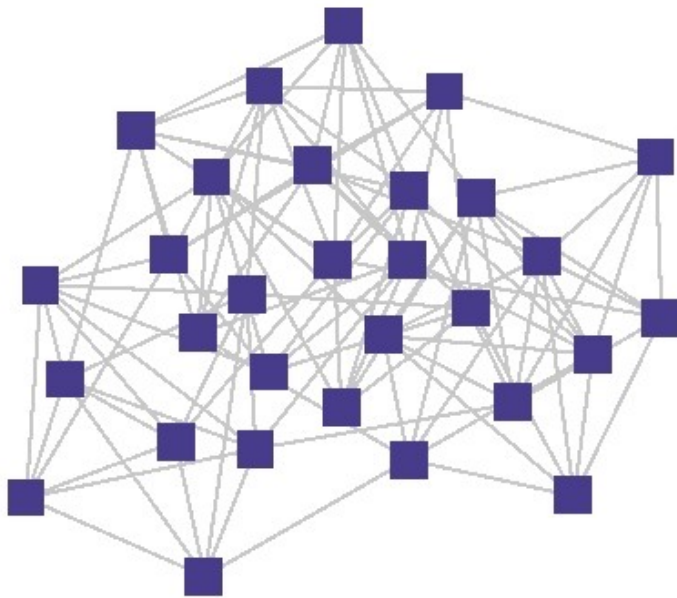
# Graph sparsification

## ✕ SPARSIFICATION

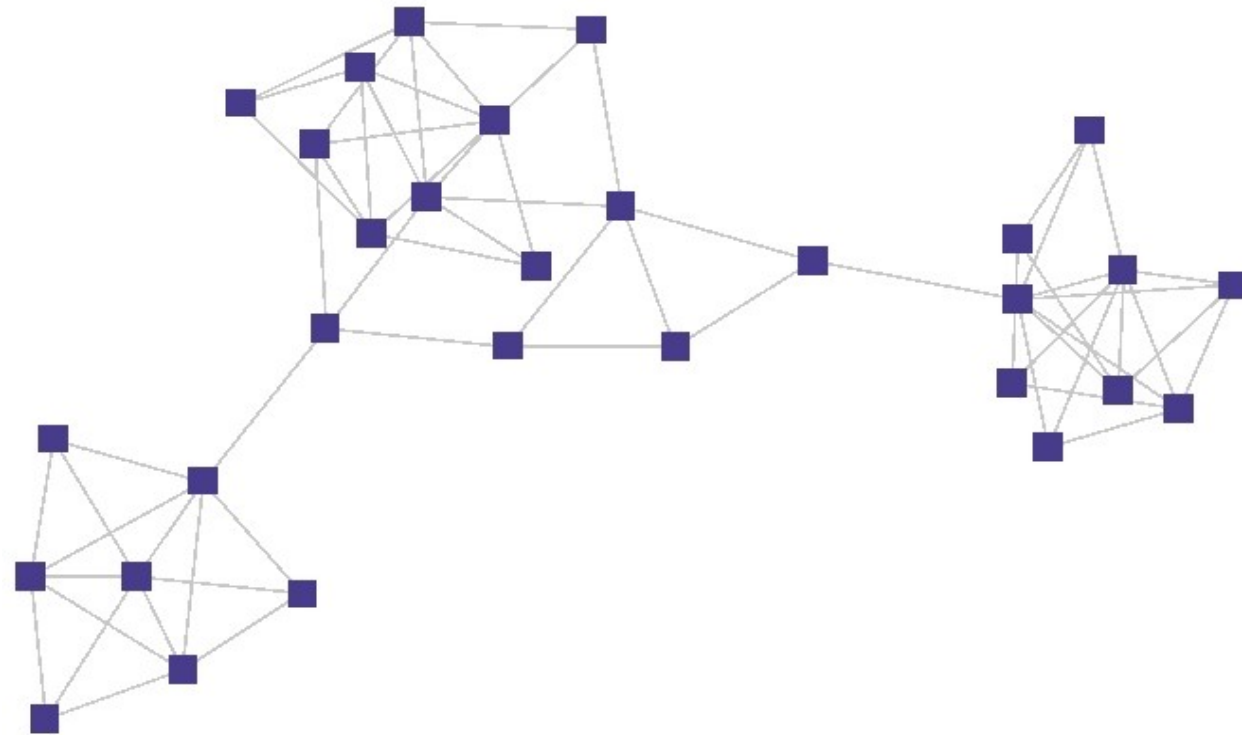
- ✕ Approximate any graph  $G$  by a sparse graph  $H$
- ✕ Nontrivial statement about  $G$
- ✕  $H$  is faster to compute with than  $G$

## ✕ IT'S NOT A TRIVIAL TASK

# Sparsification



# Original

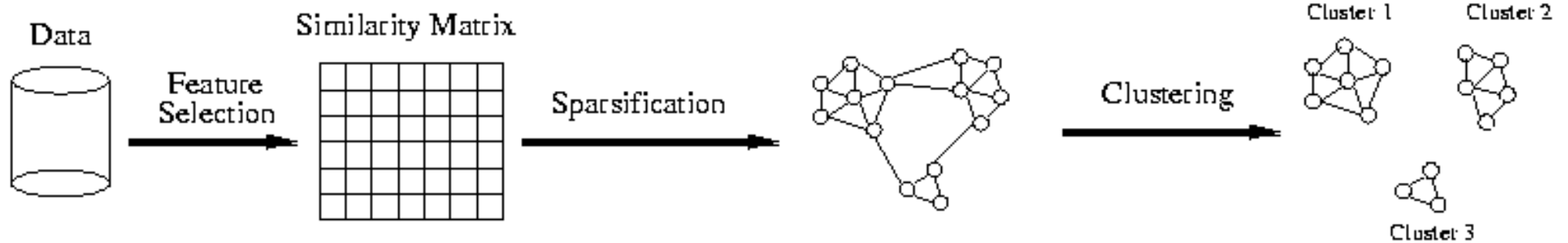


# Sparsified

# Graph-Based Clustering: Sparsification

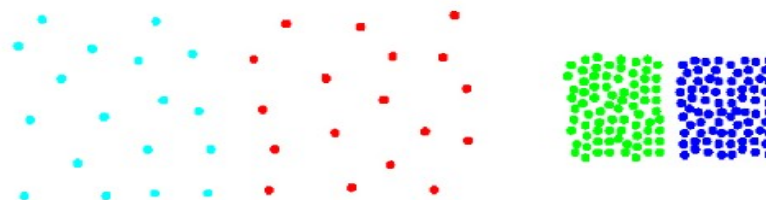
- CLUSTERING MAY WORK BETTER
  - Sparsification techniques keep the connections to the most similar (nearest) neighbors of a point while breaking the connections to less similar points.
  - The nearest neighbors of a point tend to belong to the same class as the point itself.
  - This reduces the impact of noise and outliers and sharpens the distinction between clusters.
- SPARSIFICATION FACILITATES THE USE OF GRAPH PARTITIONING ALGORITHMS (OR ALGORITHMS BASED ON GRAPH PARTITIONING ALGORITHMS).
  - Chameleon and Hypergraph-based Clustering

# Sparsification in the Clustering Process



# Chameleon: Clustering Using Dynamic Modeling

- ADAPT TO THE CHARACTERISTICS OF THE DATA SET TO FIND THE NATURAL CLUSTERS
- USE A DYNAMIC MODEL TO MEASURE THE SIMILARITY BETWEEN CLUSTERS
  - Main property is the relative closeness and relative inter-connectivity of the cluster
  - Two clusters are combined if the resulting cluster shares certain *properties* with the constituent clusters
  - The merging scheme preserves *self-similarity*
- ONE OF THE AREAS OF APPLICATION IS SPATIAL DATA

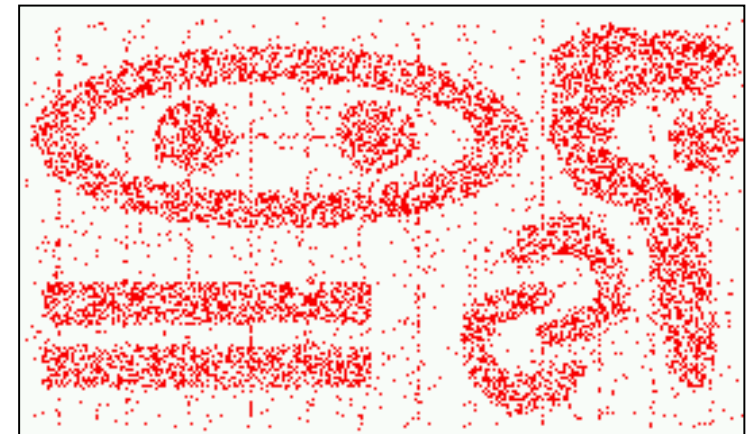
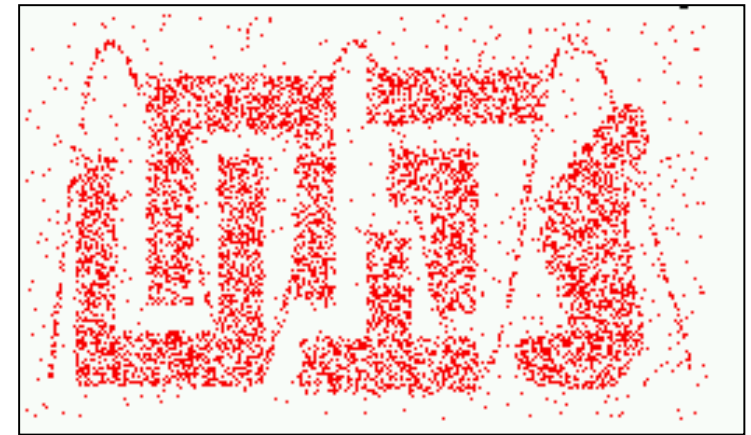




# Characteristics of Spatial Data Sets

- ✗ CLUSTERS ARE DEFINED AS DENSELY POPULATED REGIONS OF THE SPACE
- ✗ CLUSTERS HAVE ARBITRARY SHAPES, ORIENTATION, AND NON-UNIFORM SIZES
- ✗ DIFFERENCE IN DENSITIES ACROSS CLUSTERS AND VARIATION IN DENSITY WITHIN CLUSTERS
- ✗ EXISTENCE OF SPECIAL ARTIFACTS (STREAKS) AND NOISE

The clustering algorithm must address the above characteristics and also require minimal supervision.



# Chameleon: Steps

## ✗PREPROCESSING STEP:

### REPRESENT THE DATA BY A GRAPH

- ✗Given a set of points, construct the k-nearest-neighbor (k-NN) graph to capture the relationship between a point and its k nearest neighbors
- ✗Concept of neighborhood is captured dynamically (even if region is sparse)

### ✗PHASE 1: USE A MULTILEVEL GRAPH PARTITIONING ALGORITHM ON THE GRAPH TO FIND A LARGE NUMBER OF CLUSTERS OF WELL-CONNECTED VERTICES

- ✗Each cluster should contain mostly points from one “true” cluster, i.e., is a sub-cluster of a “real” cluster

# Chameleon: Steps

## x **PHASE 2**: USE HIERARCHICAL AGGLOMERATIVE CLUSTERING TO MERGE SUB-CLUSTERS

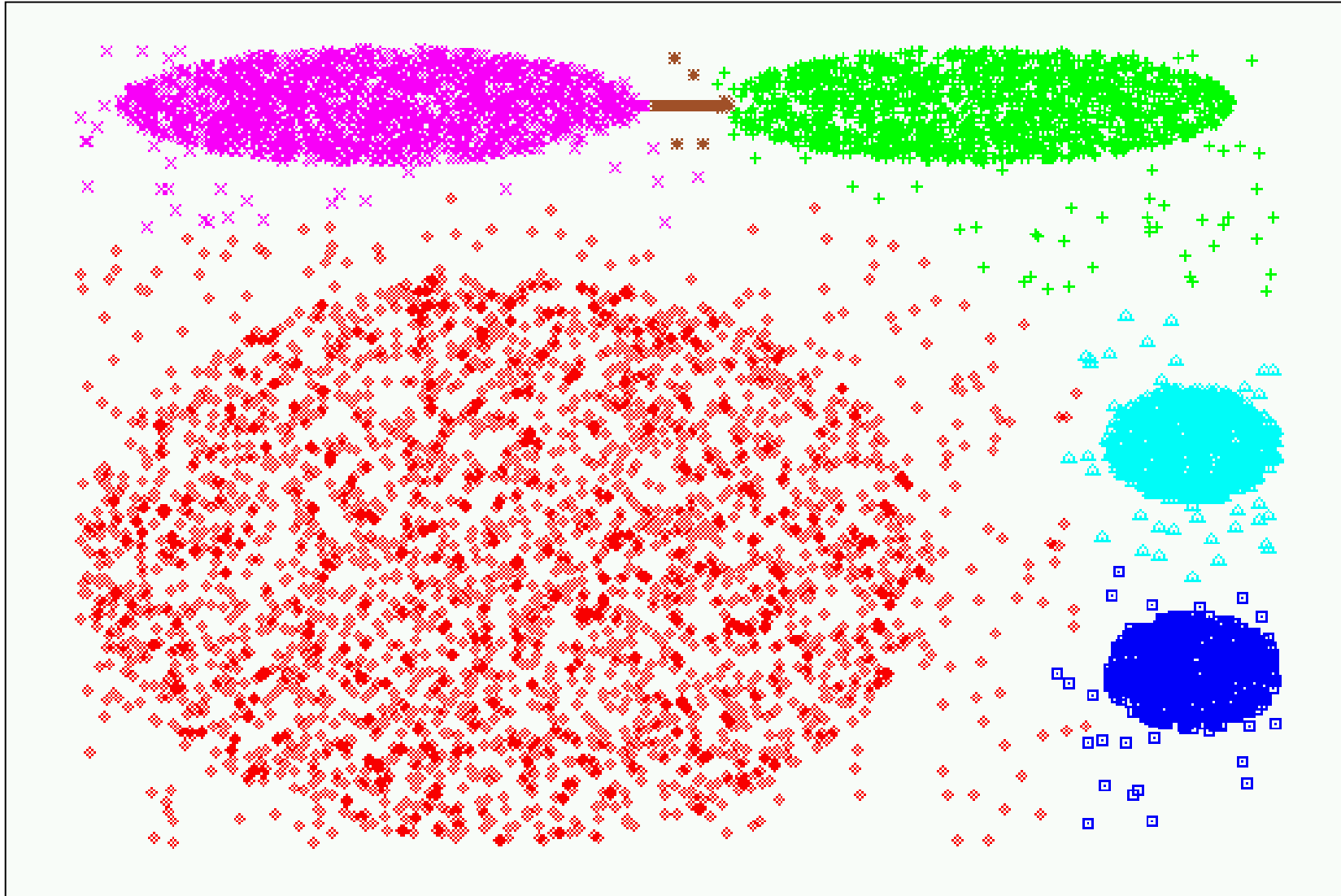
- x Two clusters are combined if the *resulting cluster shares certain properties with the constituent clusters*

- x Two key properties used to model cluster similarity:

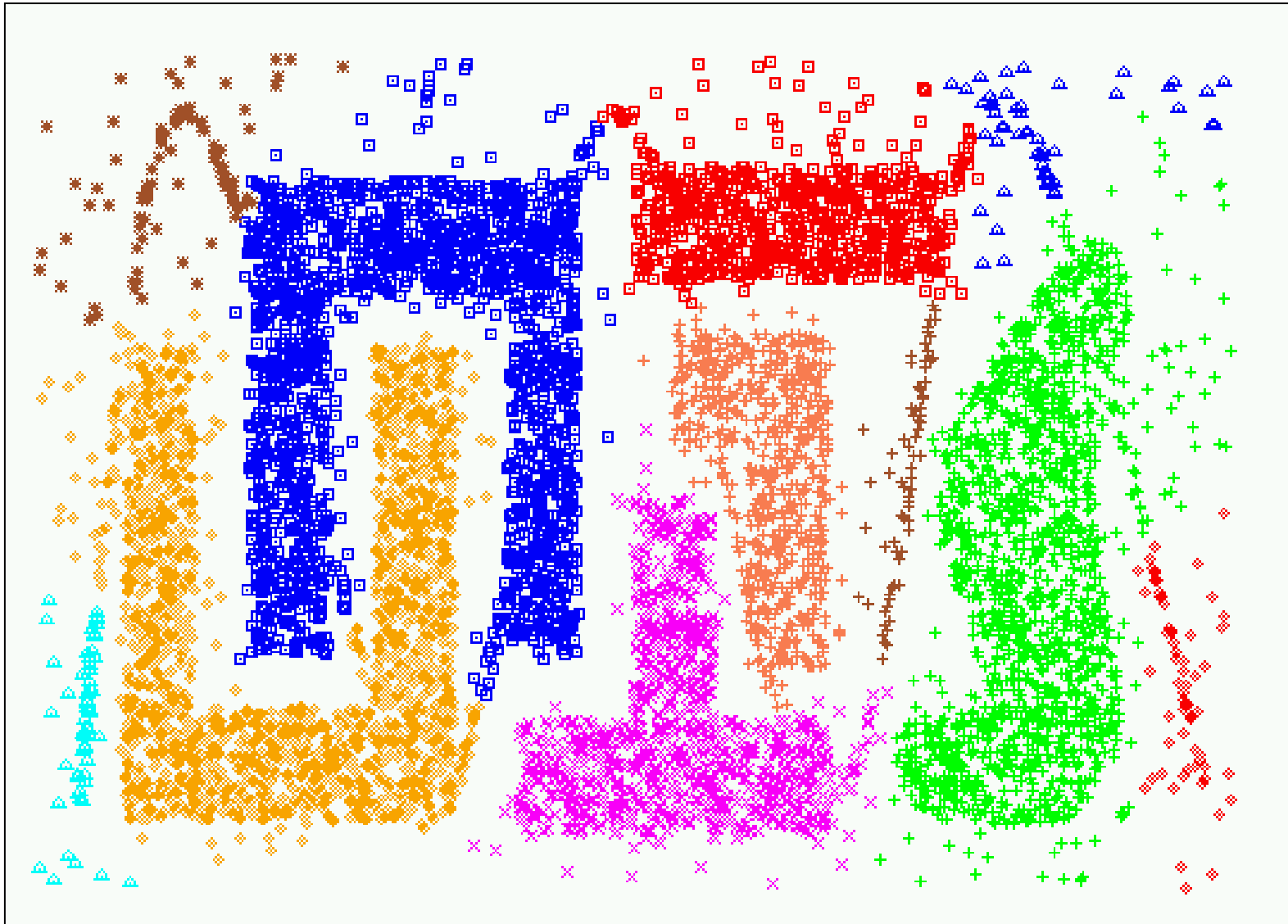
  - x **Relative Interconnectivity**: Absolute interconnectivity of two clusters normalized by the internal connectivity of the clusters

  - x **Relative Closeness**: Absolute closeness of two clusters normalized by the internal closeness of the clusters

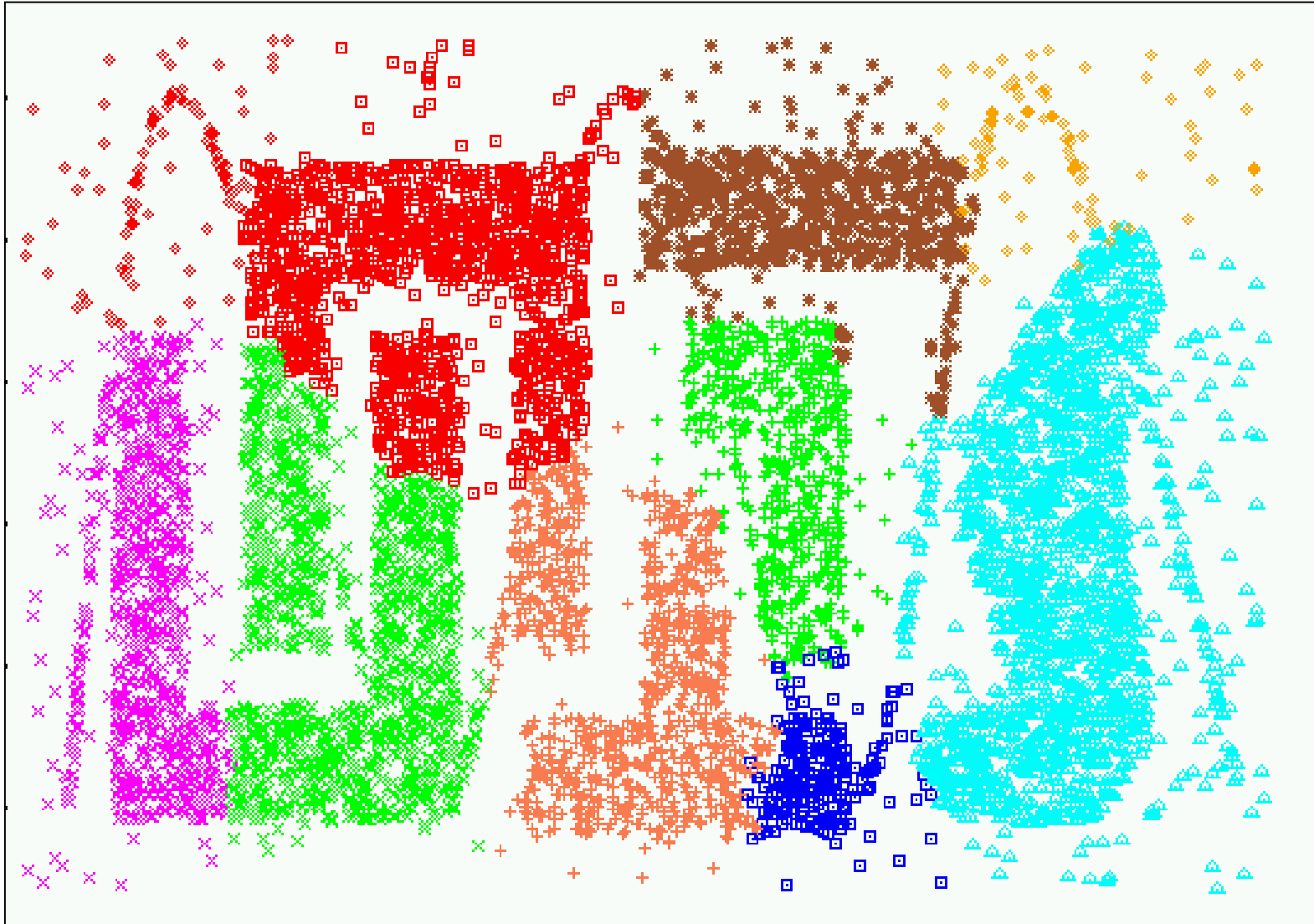
# Experimental Results: CHAMELEON



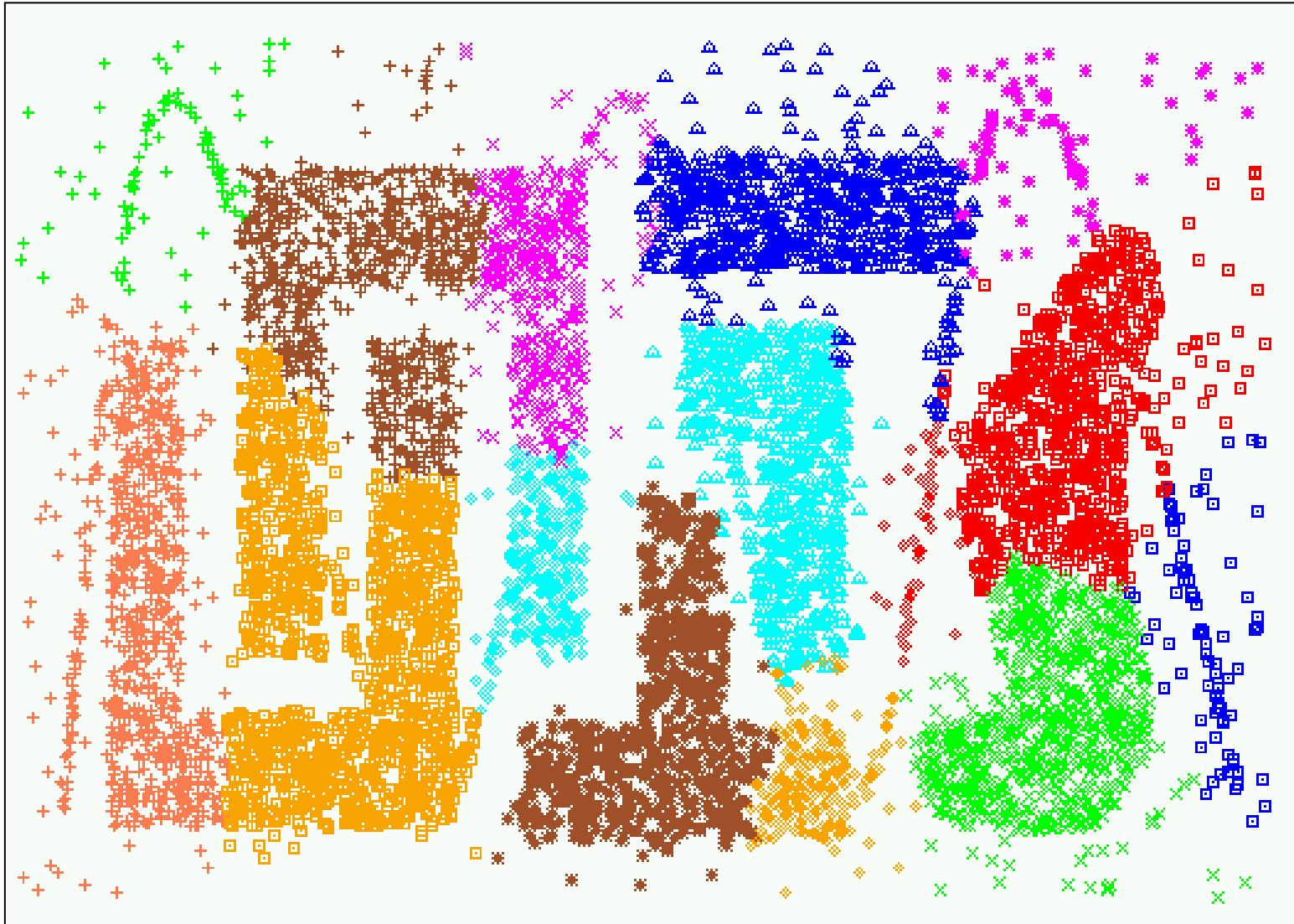
# Experimental Results: CHAMELEON



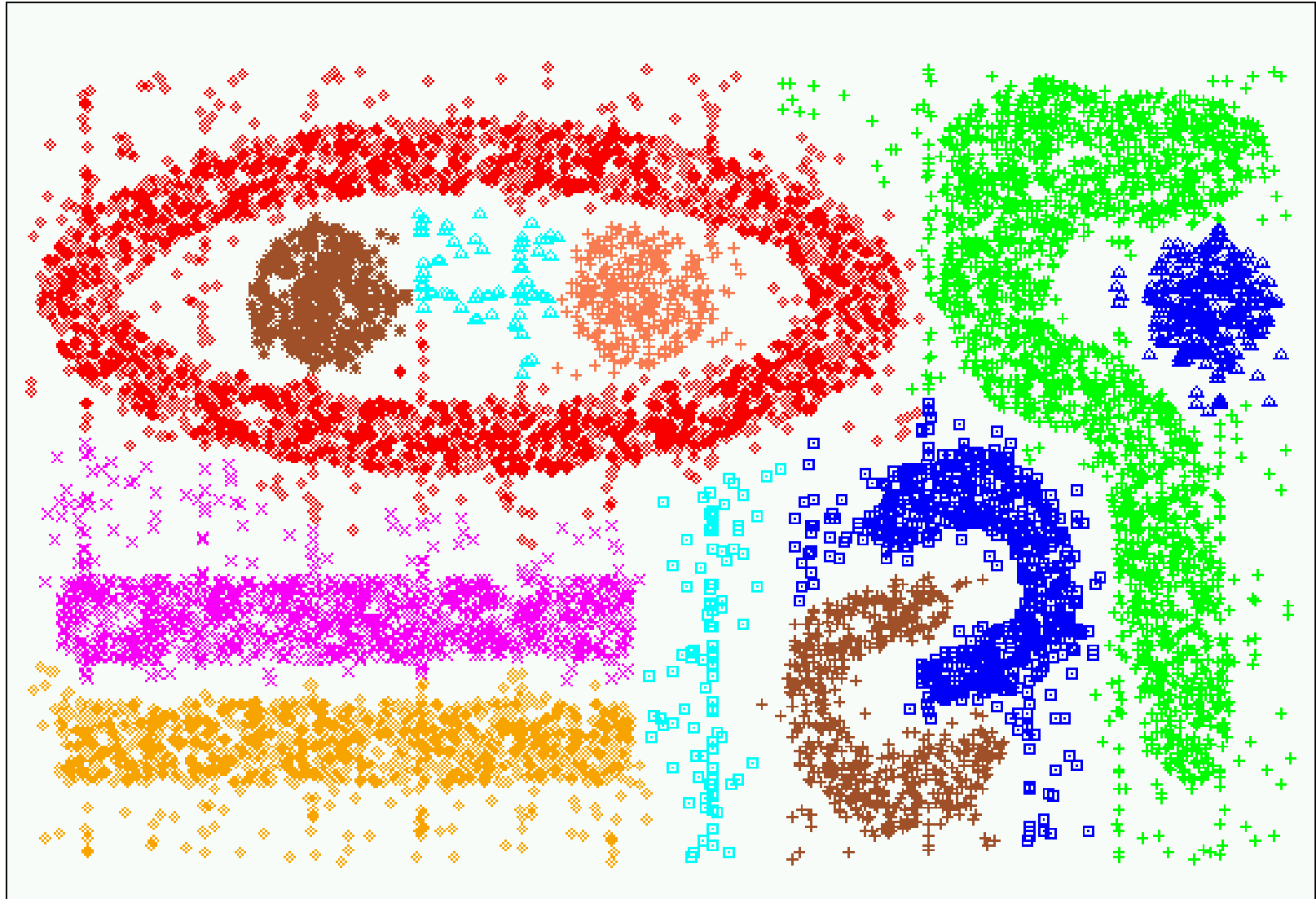
# Experimental Results: CURE (*10 clusters*)



# Experimental Results: CURE (*15 clusters*)

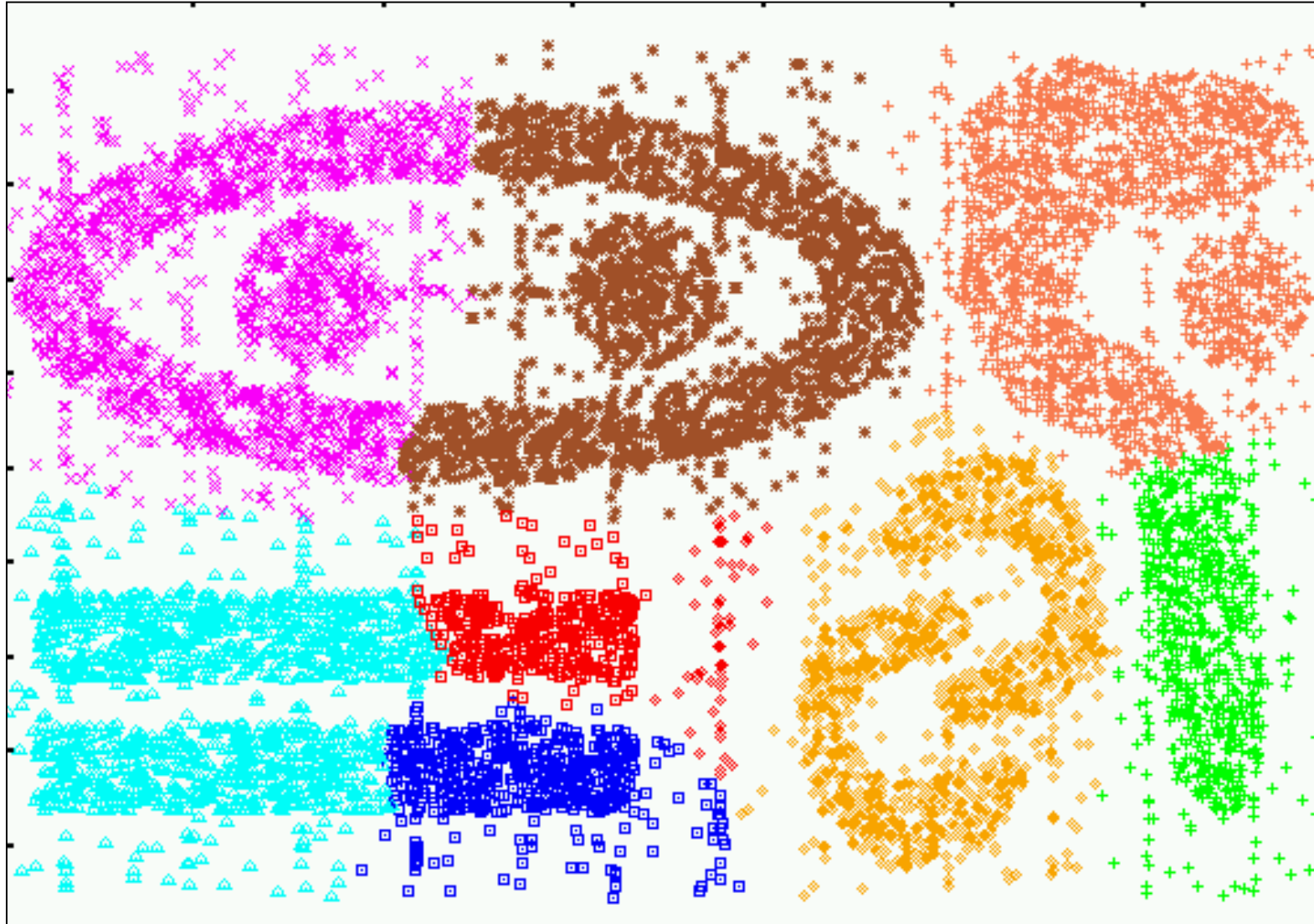


# Experimental Results: CHAMELEON

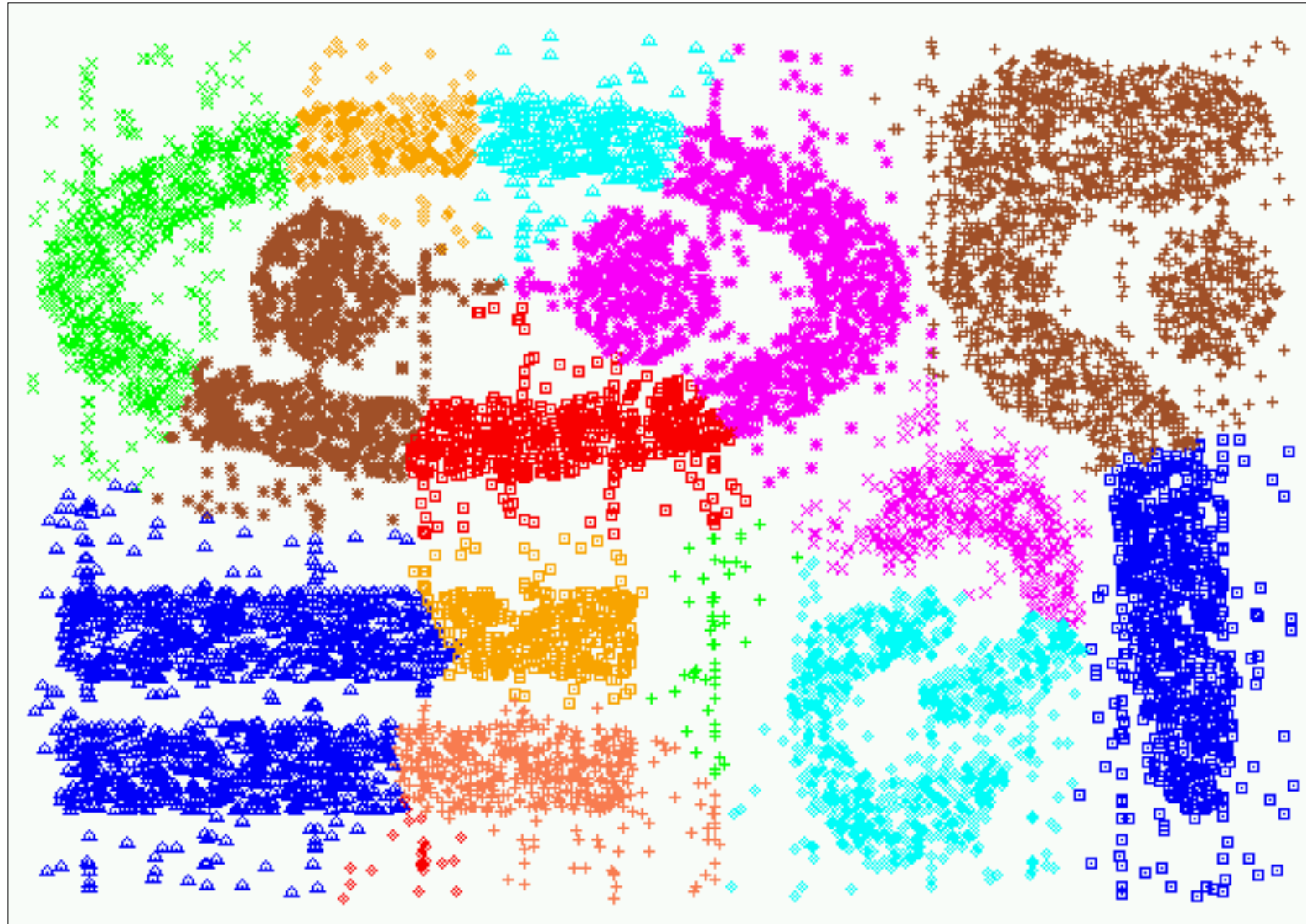




# Experimental Results: CURE (9 clusters)

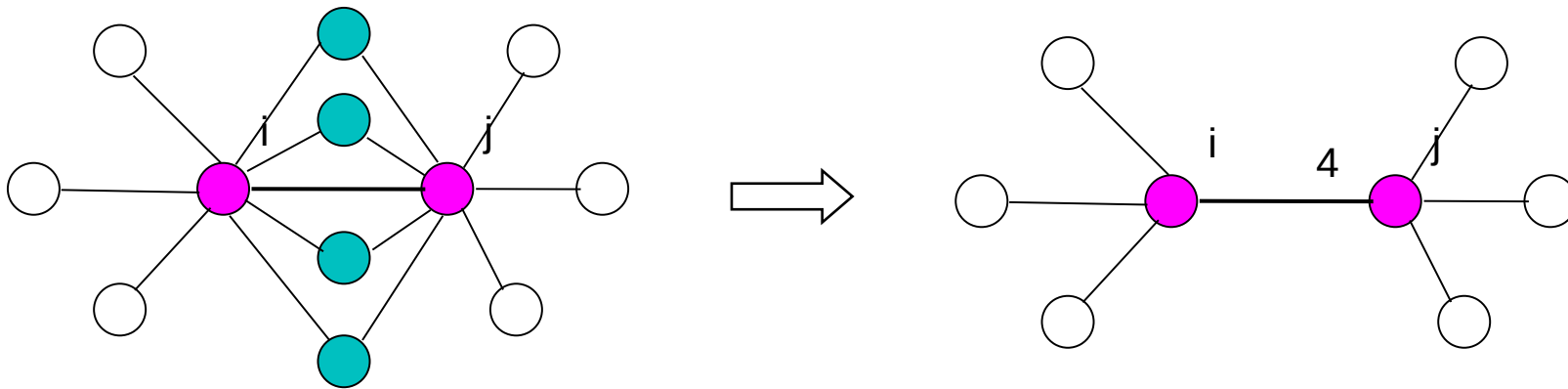


# Experimental Results: CURE (*15 clusters*)

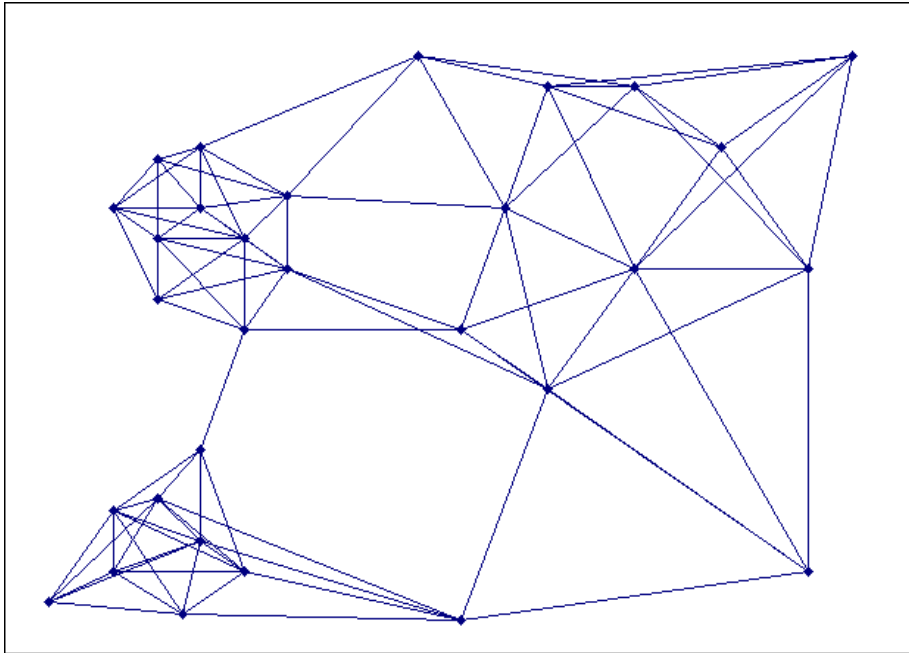


# Shared Near Neighbor Approach

**SNN graph:** the weight of an edge is the number of shared neighbors between vertices given that the vertices are connected

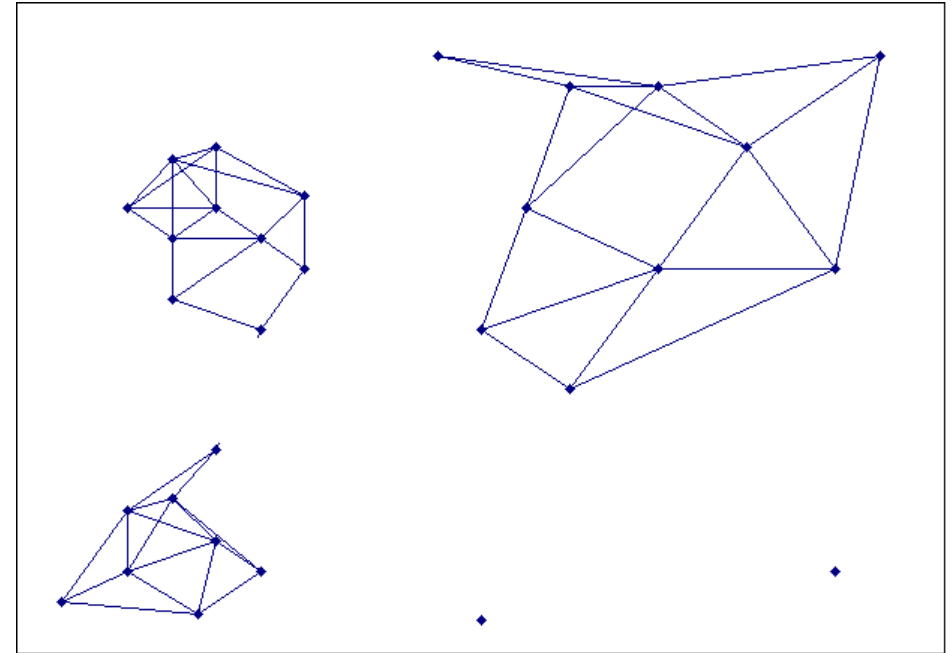


# Creating the SNN Graph



**Sparse Graph**

**Link weights are similarities  
between neighboring points**



**Shared Near Neighbor Graph**

**Link weights are number of Shared  
Nearest Neighbors**

# ROCK (RObust Clustering using linKs)

- CLUSTERING ALGORITHM FOR DATA WITH CATEGORICAL AND BOOLEAN ATTRIBUTES
  - ✗A pair of points is defined to be neighbors if their similarity is greater than some threshold
  - ✗Use a hierarchical clustering scheme to cluster the data.
- 1. OBTAIN A SAMPLE OF POINTS FROM THE DATA SET
- 2. COMPUTE THE LINK VALUE FOR EACH SET OF POINTS, I.E., TRANSFORM THE ORIGINAL SIMILARITIES (COMPUTED BY JACCARD COEFFICIENT) INTO SIMILARITIES THAT REFLECT THE NUMBER OF SHARED NEIGHBORS BETWEEN POINTS
- 3. PERFORM AN AGGLOMERATIVE HIERARCHICAL CLUSTERING ON THE DATA USING THE “NUMBER OF SHARED NEIGHBORS” AS SIMILARITY MEASURE AND MAXIMIZING “THE SHARED NEIGHBORS” OBJECTIVE FUNCTION
- 4. ASSIGN THE REMAINING POINTS TO THE CLUSTERS THAT HAVE BEEN FOUND

# Density-based Clustering

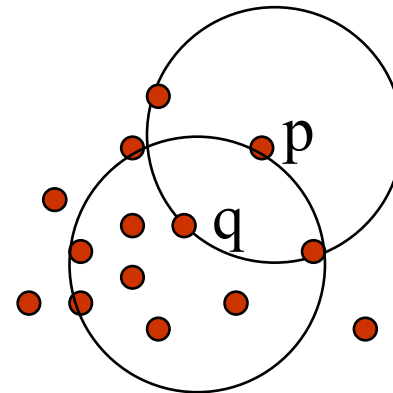
# Density-Based Clustering Methods

- CLUSTERING BASED ON DENSITY (LOCAL CLUSTER CRITERION), SUCH AS DENSITY-CONNECTED POINTS
- MAJOR FEATURES:
  - Discover clusters of arbitrary shape
  - Handle noise
  - One scan
  - Need density parameters as termination condition
- 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)

# Density-Based Clustering: Basic Concepts

- TWO PARAMETERS:
  - *Eps*: Maximum radius of the neighbourhood
  - *MinPts*: Minimum number of points in an Eps-neighbourhood of that point
- $N_{Eps}(p)$ :  $\{Q \text{ BELONGS TO } D \mid DIST(p, Q) \leq Eps\}$
- **DIRECTLY DENSITY-REACHABLE**: A POINT  $p$  IS DIRECTLY DENSITY-REACHABLE FROM A POINT  $q$  W.R.T.  $Eps$ ,  $MINPTS$  IF
  - $p$  belongs to  $N_{Eps}(q)$
  - core point condition:

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



MinPts = 5

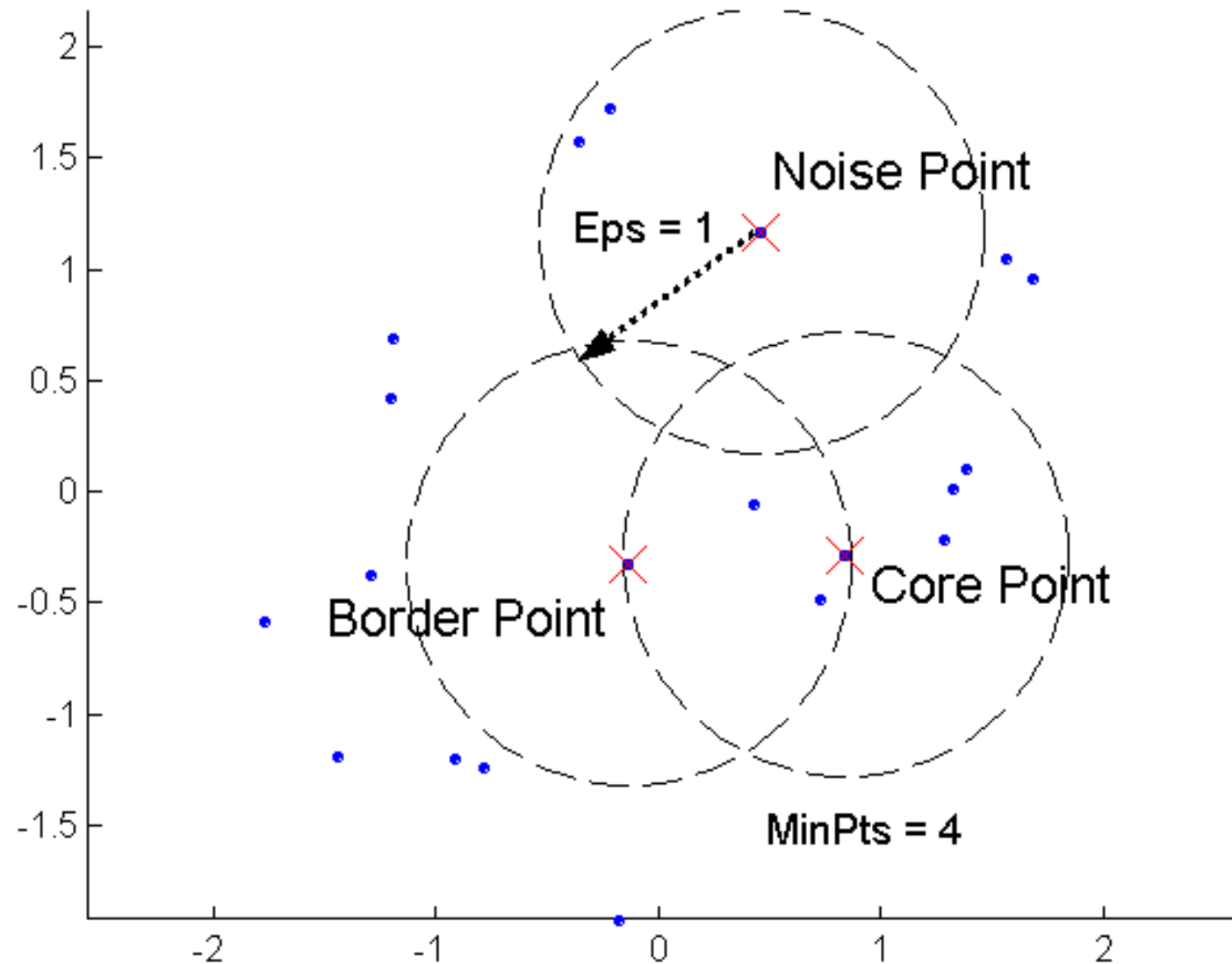
Eps = 1 cm



# Basic concepts

- ✕ 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

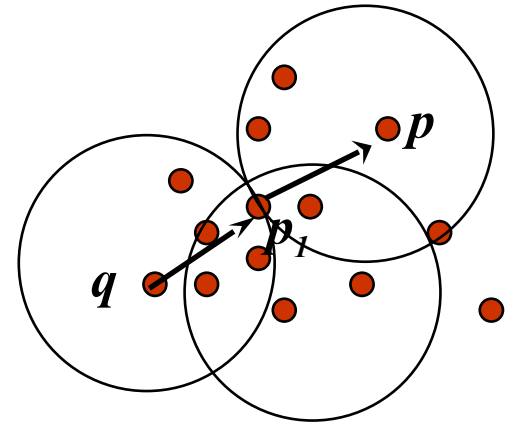
# Core, Border, and Noise Points



# Density-Reachable and Density-Connected

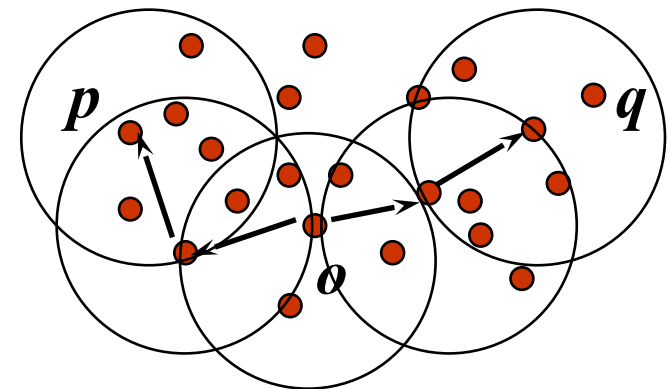
## ■ DENSITY-REACHABLE:

- A point  $p$  is **density-reachable** from a point  $q$  w.r.t.  $Eps$ ,  $MinPts$  if there is a chain of points  $p_1, \dots, p_n$ ,  $p_1 = q$ ,  $p_n = p$  such that  $p_{i+1}$  is directly density-reachable from  $p_i$



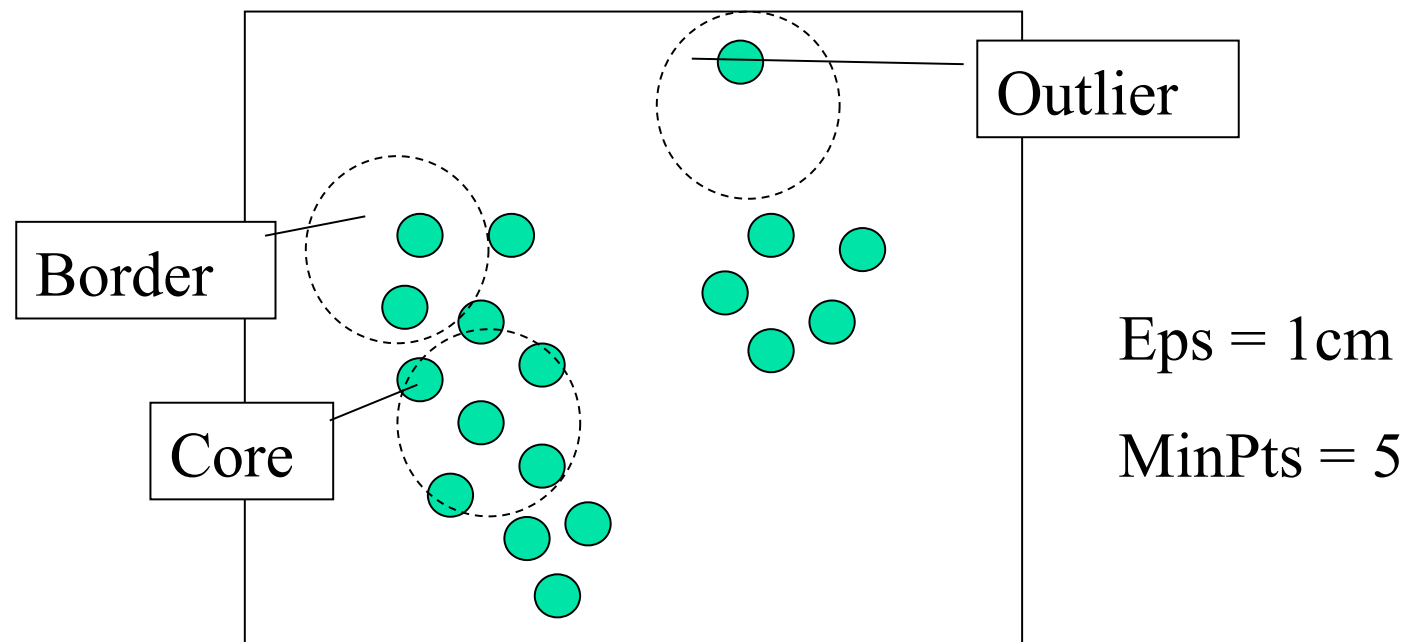
## ■ DENSITY-CONNECTED

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



# DBSCAN: Density Based Spatial Clustering of Applications with Noise

- RELIES ON A *DENSITY-BASED* NOTION OF CLUSTER: A *CLUSTER* IS DEFINED AS A MAXIMAL SET OF DENSITY-CONNECTED POINTS
- DISCOVERS CLUSTERS OF ARBITRARY SHAPE IN SPATIAL DATABASES WITH NOISE



# DBSCAN: The Algorithm

- ARBITRARY SELECT A POINT  $P$
- RETRIEVE ALL POINTS DENSITY-REACHABLE FROM  $P$  W.R.T.  $Eps$  AND  $MinPts$ .
- IF  $P$  IS A CORE POINT, A CLUSTER IS FORMED.
- IF  $P$  IS A BORDER POINT, NO POINTS ARE DENSITY-REACHABLE FROM  $P$  AND DBSCAN VISITS THE NEXT POINT OF THE DATABASE.
- CONTINUE THE PROCESS UNTIL ALL OF THE POINTS HAVE BEEN PROCESSED.

# DBSCAN Algorithm

- ✗ ELIMINATE NOISE POINTS
- ✗ PERFORM CLUSTERING ON THE REMAINING POINTS

$current\_cluster\_label \leftarrow 1$

**for** all core points **do**

**if** the core point has no cluster label **then**

$current\_cluster\_label \leftarrow current\_cluster\_label + 1$

        Label the current core point with cluster label  $current\_cluster\_label$

**end if**

**for** all points in the  $Eps$ -neighborhood, except  $i^{th}$  the point itself **do**

**if** the point does not have a cluster label **then**

            Label the point with cluster label  $current\_cluster\_label$

**end if**

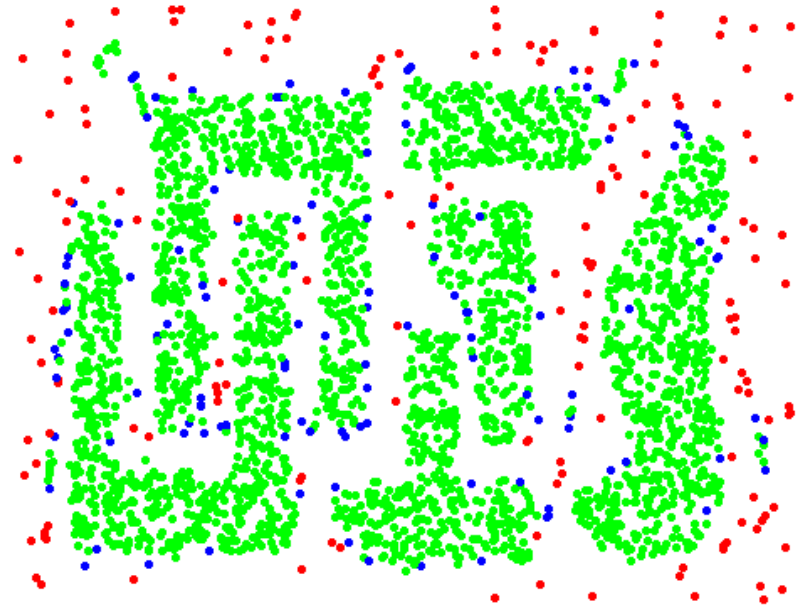
**end for**

**end for**

# DBSCAN: Core, Border and Noise Points



ORIGINAL POINTS



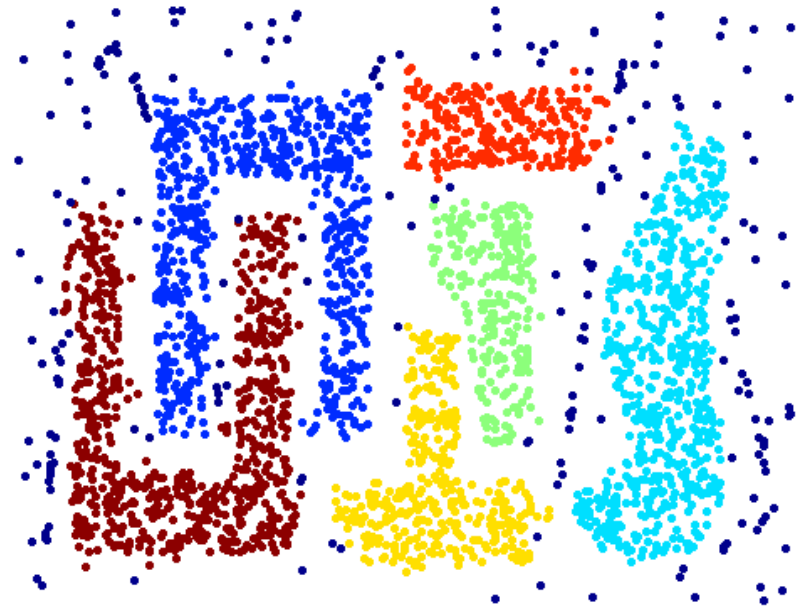
POINT TYPES: CORE, BORDER  
AND NOISE

EPS = 10, MINPTS = 4

# When DBSCAN Works Well



ORIGINAL POINTS



CLUSTERS

- RESISTANT TO NOISE
- CAN HANDLE CLUSTERS OF DIFFERENT SHAPES AND SIZES



# DBSCAN: Sensitive to Parameters

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

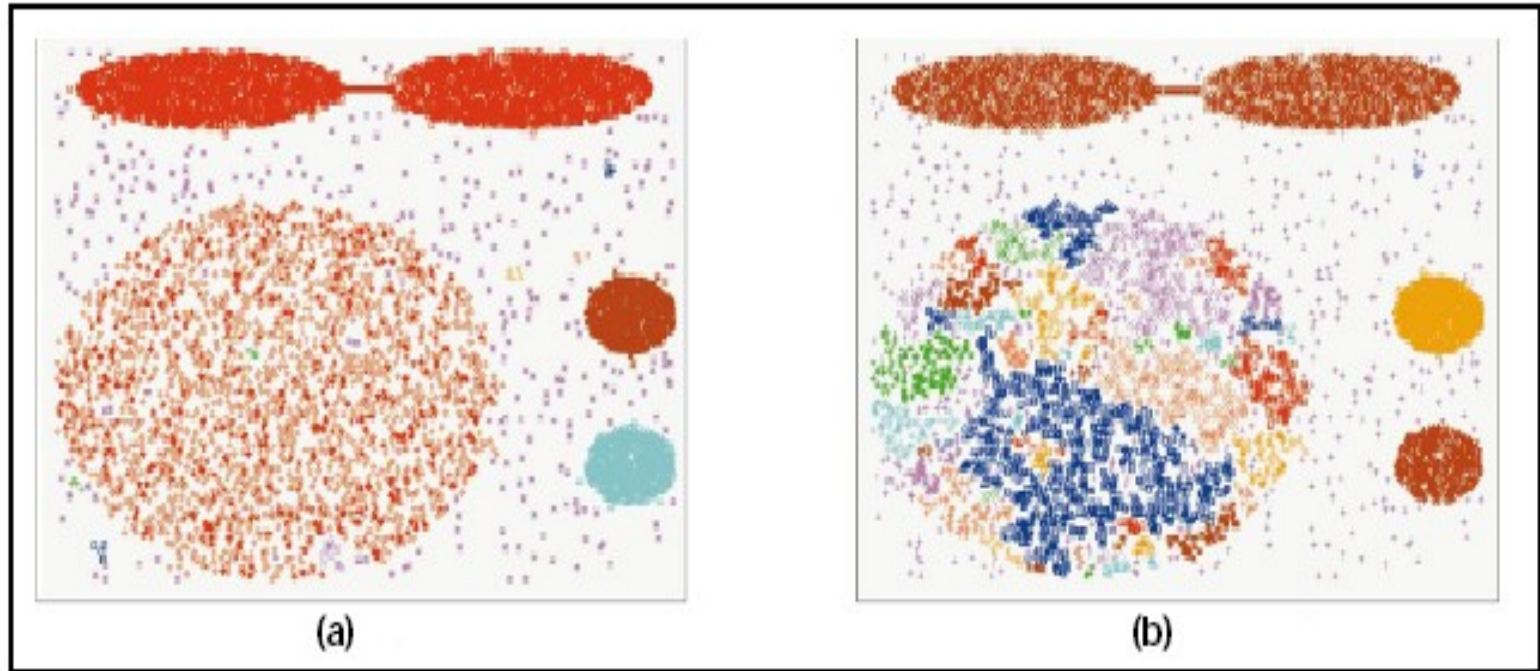
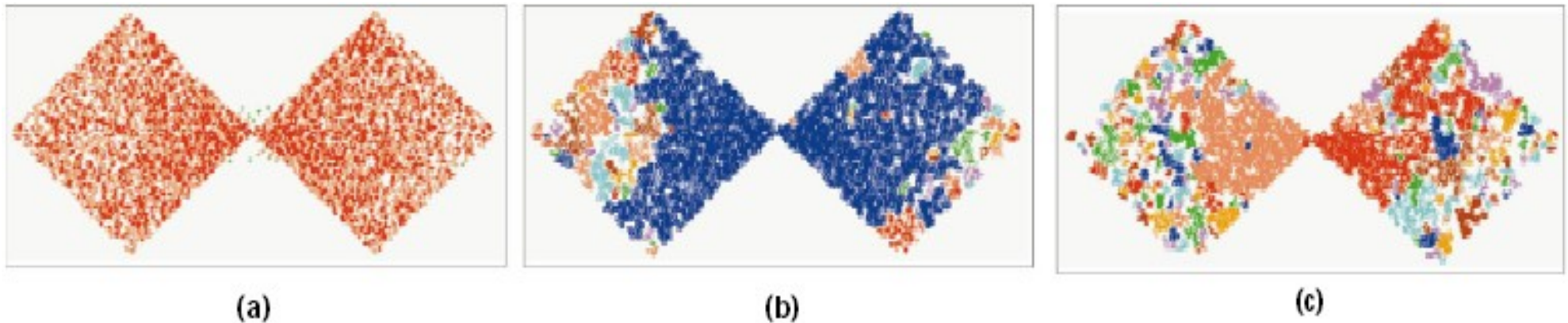
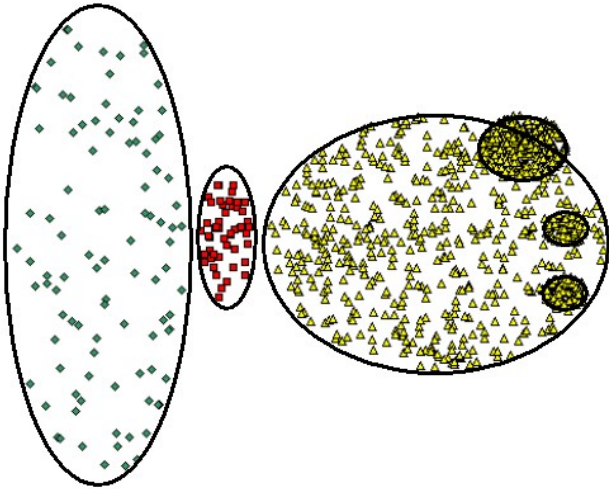


Figure 9. DBScan results for DS2 with MinPts at 4 and Eps at (a) 5.0, (b) 3.5, and (c) 3.0.

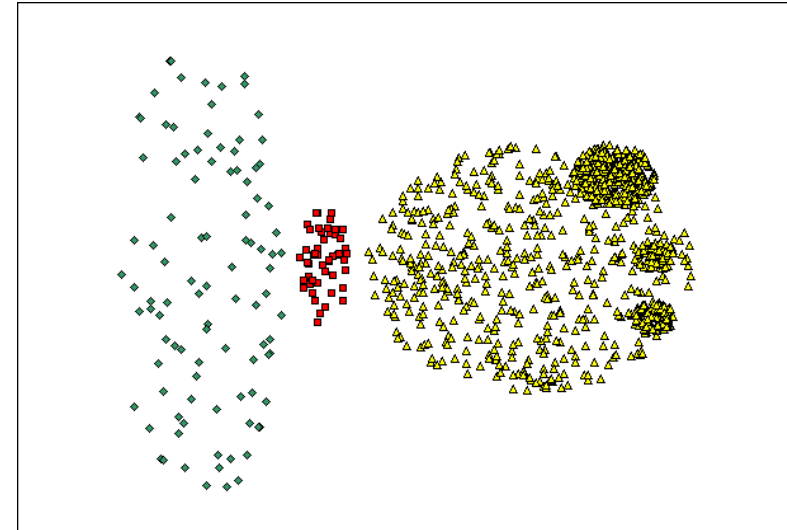


# When DBSCAN Does NOT Work Well

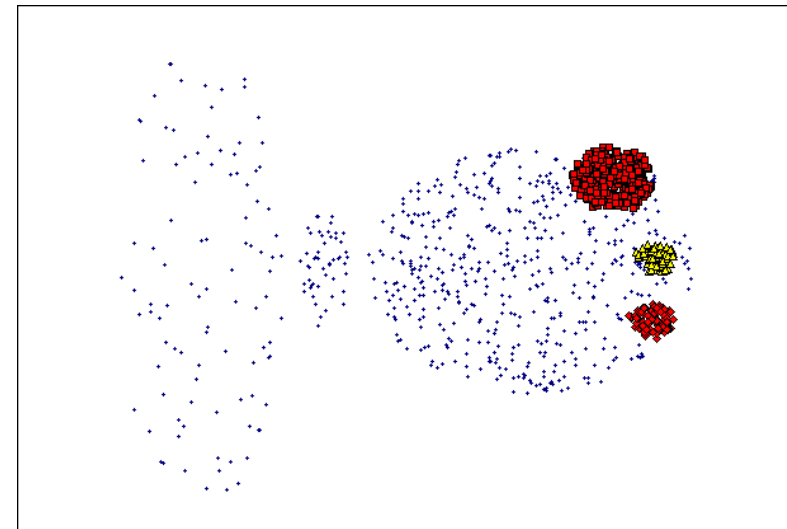


ORIGINAL POINTS

- VARYING DENSITIES
- HIGH-DIMENSIONAL DATA



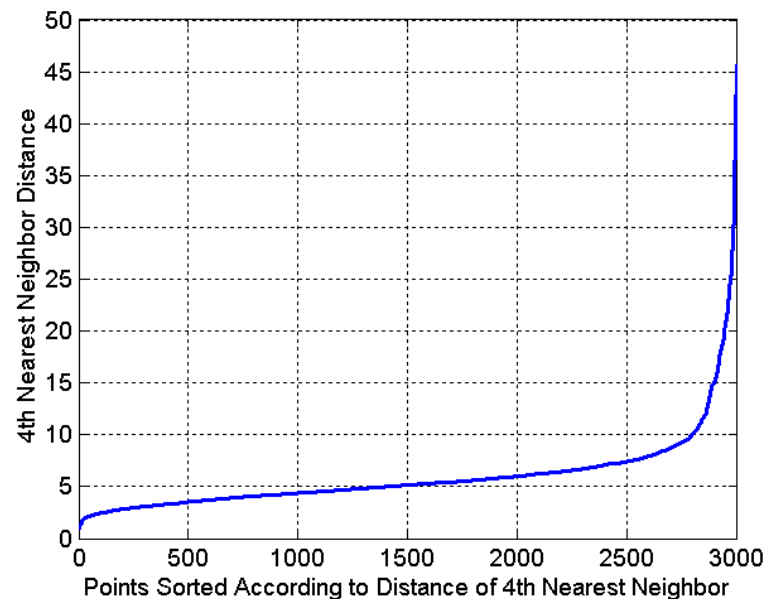
(MinPts=4, Eps=9.75).



(MinPts=4, Eps=9.92)

# DBSCAN: Determining EPS and MinPts

- ✗ IDEA IS THAT FOR POINTS IN A CLUSTER, THEIR K-TH NEAREST NEIGHBORS ARE AT ROUGHLY THE SAME DISTANCE
- ✗ NOISE POINTS HAVE THE K-TH NEAREST NEIGHBOR AT FARTHER DISTANCE
- ✗ SO, PLOT SORTED DISTANCE OF EVERY POINT TO ITS K-TH NEAREST NEIGHBOR



# OPTICS: A Cluster-Ordering Method (1999)

- OPTICS: ORDERING POINTS TO IDENTIFY THE CLUSTERING STRUCTURE
  - Ankerst, Breunig, Kriegel, and Sander (SIGMOD'99)
  - Produces a special order of the database wrt its density-based clustering structure
  - This cluster-ordering contains info equiv to the density-based clusterings corresponding to a broad range of parameter settings
  - Good for both automatic and interactive cluster analysis, including finding intrinsic clustering structure
  - Can be represented graphically or using visualization techniques

# DENCLUE: Using Statistical Density Functions

- DENSITY-BASED CLUSTERING BY HINNEBURG & KEIM (KDD'98)

- USING STATISTICAL DENSITY FUNCTIONS:  $f_{Gaussian}(x, y) = e^{-\frac{d(x, y)^2}{2\sigma^2}}$

$$f_{Gaussian}^D(x) = \sum_{i=1}^N e^{-\frac{d(x, x_i)^2}{2\sigma^2}}$$

- MAJOR FEATURES

- Solid mathematical foundation  $\nabla f_{Gaussian}^D(x, x_i) = \sum_{i=1}^N (x_i - x) \cdot e^{-\frac{d(x, x_i)^2}{2\sigma^2}}$
- Good for data sets with large amounts of noise
- Allows a compact mathematical description of arbitrarily shaped clusters in high-dimensional data sets
- Significant faster than existing algorithm (e.g., DBSCAN)
- But needs a large number of parameters

# Denclue: Technical Essence

- USES GRID CELLS BUT ONLY KEEPS INFORMATION ABOUT GRID CELLS THAT DO ACTUALLY CONTAIN DATA POINTS AND MANAGES THESE CELLS IN A TREE-BASED ACCESS STRUCTURE
- INFLUENCE FUNCTION: DESCRIBES THE IMPACT OF A DATA POINT WITHIN ITS NEIGHBORHOOD
- OVERALL DENSITY OF THE DATA SPACE CAN BE CALCULATED AS THE SUM OF THE INFLUENCE FUNCTION OF ALL DATA POINTS
- CLUSTERS CAN BE DETERMINED MATHEMATICALLY BY IDENTIFYING DENSITY ATTRACTORS
- DENSITY ATTRACTORS ARE LOCAL MAXIMAL OF THE OVERALL DENSITY FUNCTION

# Clustering High-Dimensional Data

## xCLUSTERING HIGH-DIMENSIONAL DATA

- xMany applications: text documents, DNA micro-array data
- xMajor challenges:
  - xMany irrelevant dimensions may mask clusters
  - xDistance measure becomes meaningless—due to equi-distance
  - xClusters may exist only in some subspaces

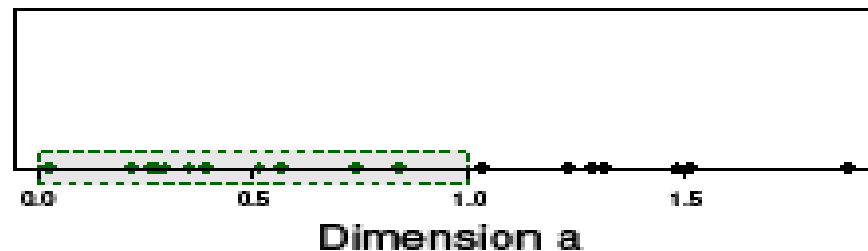
## xMETHODS

- xFeature transformation: only effective if most dimensions are relevant
  - xPCA & SVD useful only when features are highly correlated/redundant
- xFeature selection: wrapper or filter approaches
  - xuseful to find a subspace where the data have nice clusters
- xSubspace-clustering: find clusters in all the possible subspaces
  - xCLIQUE, ProClus, and frequent pattern-based clustering

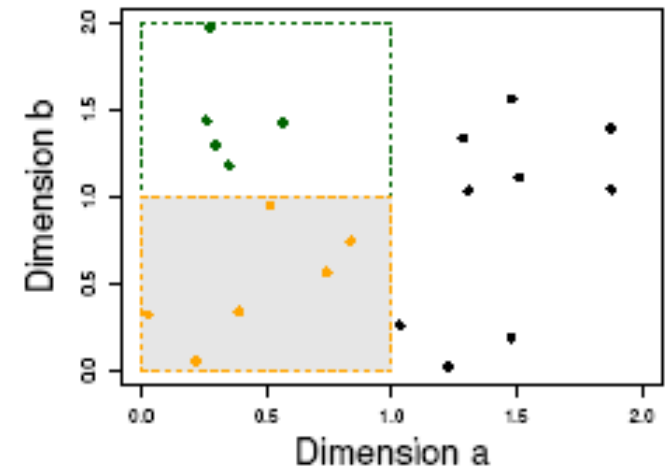


# The Curse of Dimensionality

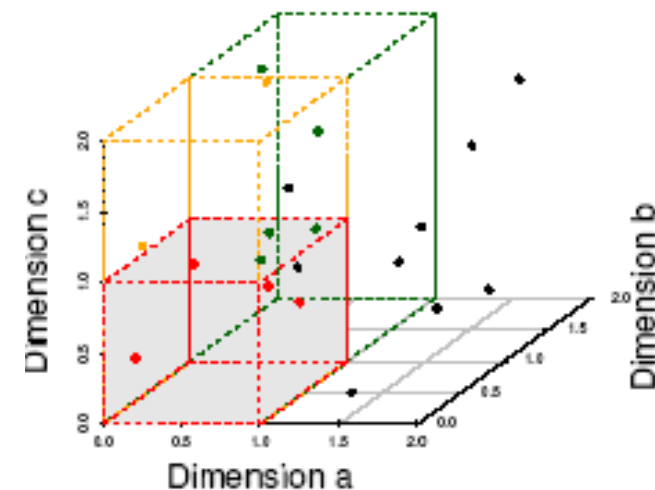
- DATA IN ONLY ONE DIMENSION IS RELATIVELY PACKED
- ADDING A DIMENSION “STRETCH” THE POINTS ACROSS THAT DIMENSION, MAKING THEM FURTHER APART
- ADDING MORE DIMENSIONS WILL MAKE THE POINTS FURTHER APART—HIGH DIMENSIONAL DATA IS EXTREMELY SPARSE
- DISTANCE MEASURE BECOMES MEANINGLESS—DUE TO EQUI-DISTANCE



(a) 11 Objects in One Unit Bin



(b) 6 Objects in One Unit Bin

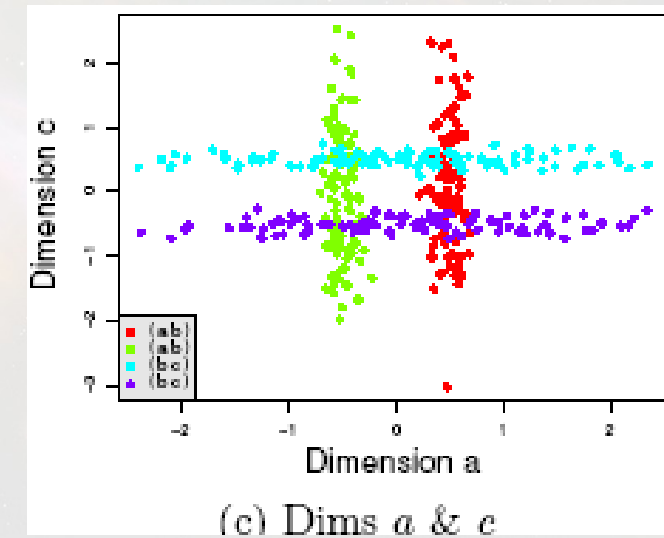
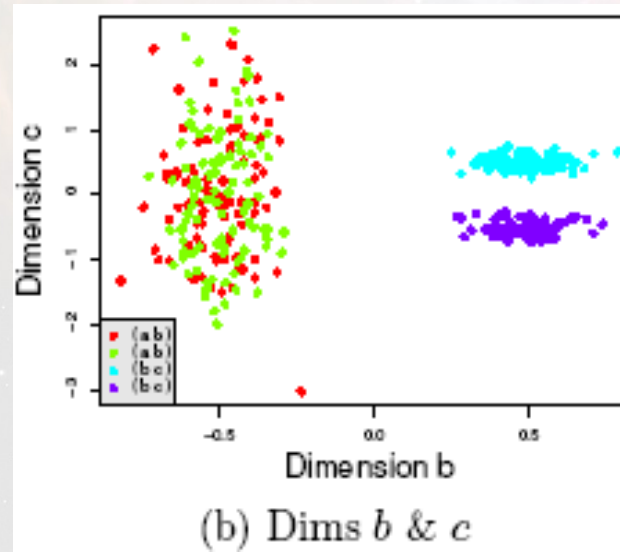
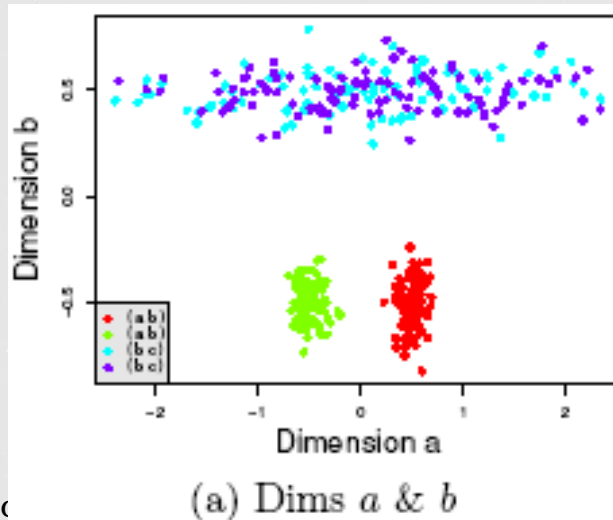
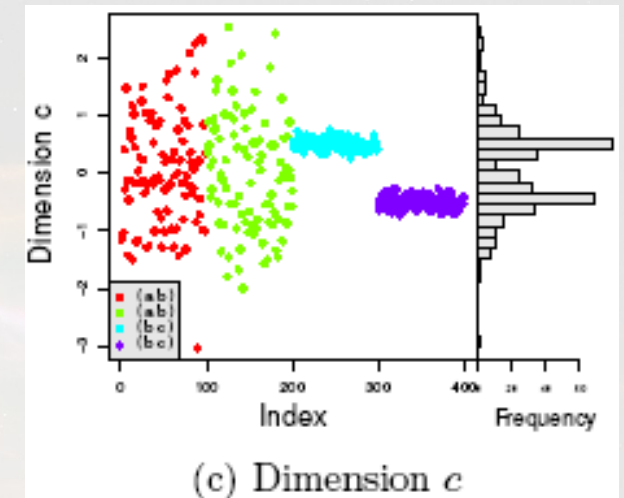
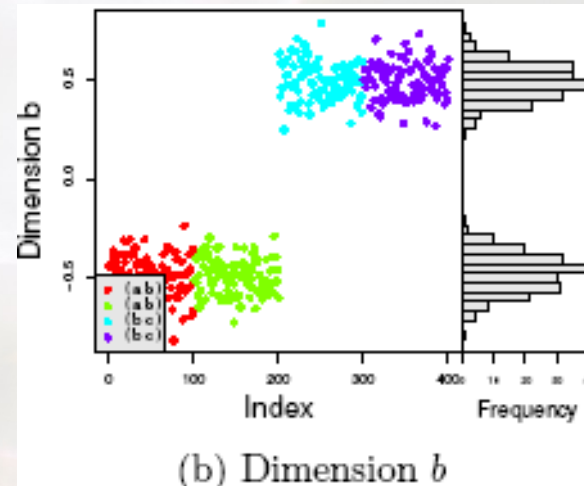
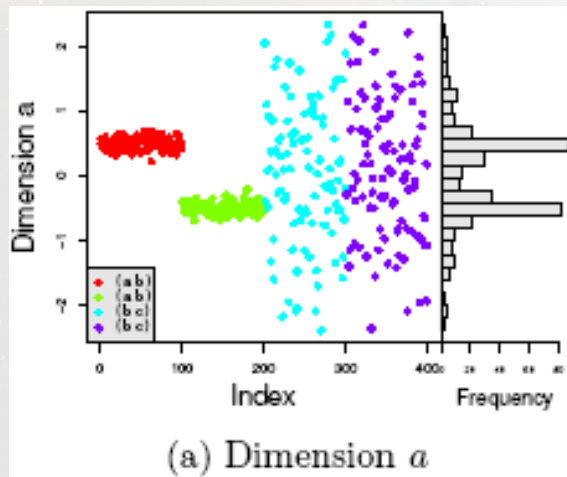
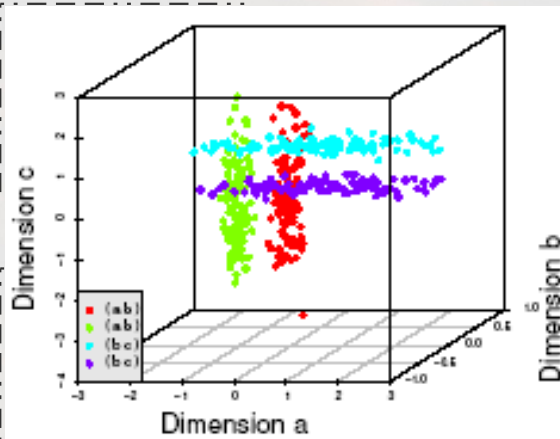


(c) 4 Objects in One Unit Bin



# Why Subspace Clustering?

- CLUSTERS MAY EXIST ONLY IN SOME SUBSPACES
- SUBSPACE-CLUSTERING: FIND CLUSTERS IN ALL THE SUBSPACES

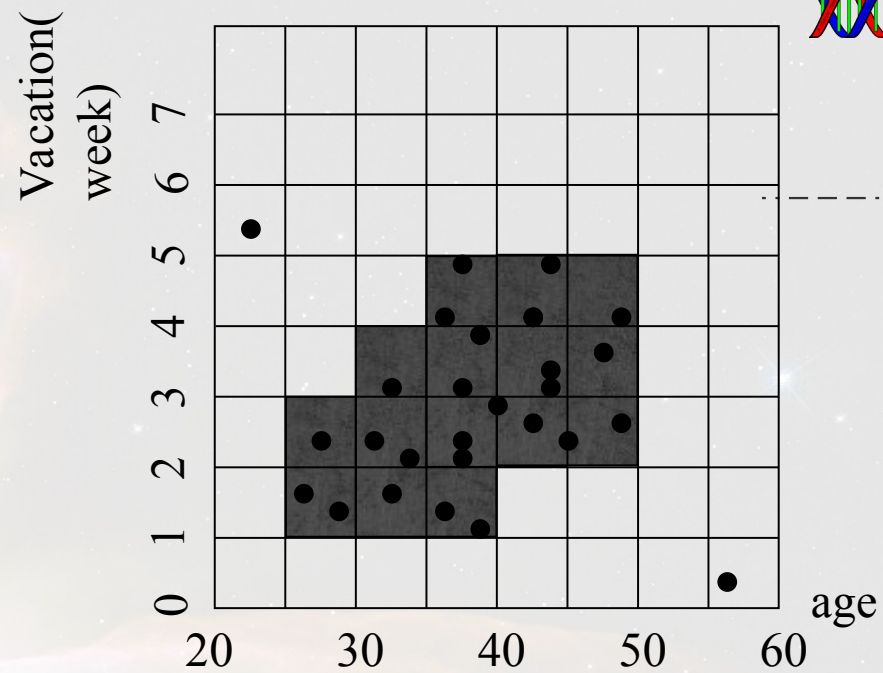
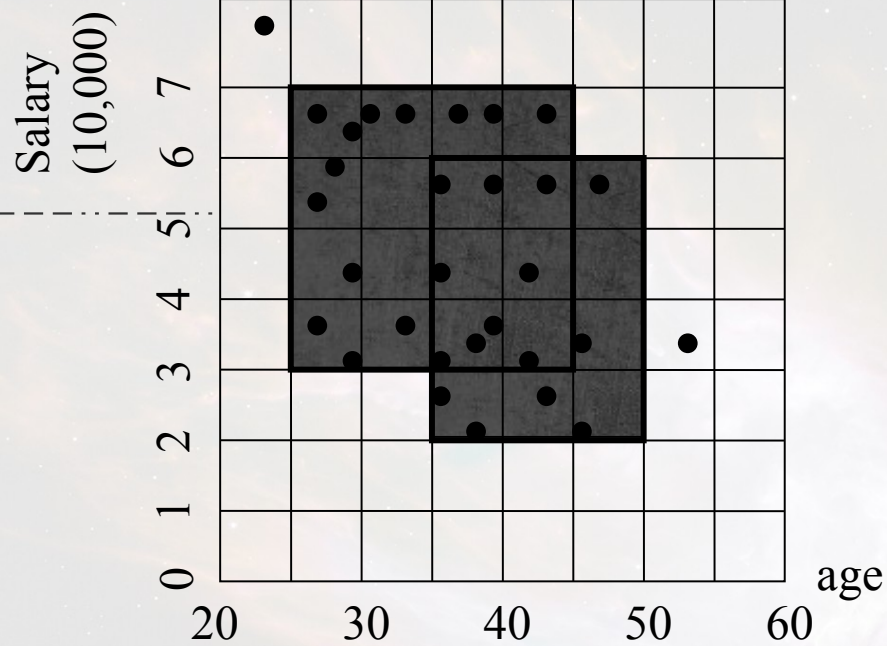


# CLIQUE (Clustering In QUES)

- AGRAWAL, GEHRKE, GUNOPULOS, RAGHAVAN (SIGMOD'98)
- AUTOMATICALLY IDENTIFYING SUBSPACES OF A HIGH DIMENSIONAL DATA SPACE THAT ALLOW BETTER CLUSTERING THAN ORIGINAL SPACE
- CLIQUE CAN BE CONSIDERED AS BOTH DENSITY-BASED AND GRID-BASED
  - It partitions each dimension into the same number of equal length interval
  - It partitions an m-dimensional data space into non-overlapping rectangular units
  - A unit is dense if the fraction of total data points contained in the unit exceeds the input model parameter
  - A cluster is a maximal set of connected dense units within a subspace

# CLIQUE: The Major Steps

- PARTITION THE DATA SPACE AND FIND THE NUMBER OF POINTS THAT LIE INSIDE EACH CELL OF THE PARTITION.
- IDENTIFY THE SUBSPACES THAT CONTAIN CLUSTERS USING THE APRIORI PRINCIPLE
- IDENTIFY CLUSTERS
  - Determine dense units in all subspaces of interests
  - Determine connected dense units in all subspaces of interests.
- GENERATE MINIMAL DESCRIPTION FOR THE CLUSTERS
  - Determine maximal regions that cover a cluster of connected dense units for each cluster
  - Determination of minimal cover for each cluster



$\tau = 3$



# Strength and Weakness of *CLIQUE*

## ■ STRENGTH

- *automatically* finds subspaces of the highest dimensionality such that high density clusters exist in those subspaces
- *insensitive* to the order of records in input and does not presume some canonical data distribution
- scales *linearly* with the size of input and has good scalability as the number of dimensions in the data increases

## ■ WEAKNESS

- The accuracy of the clustering result may be degraded at the expense of simplicity of the method

# Frequent Pattern-Based Approach

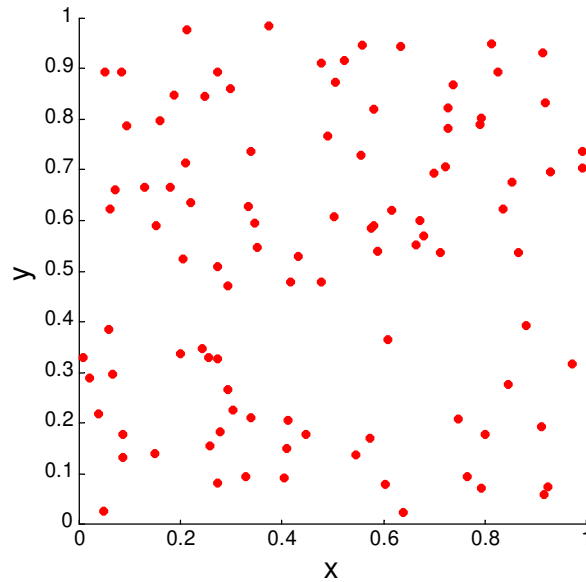
- CLUSTERING HIGH-DIMENSIONAL SPACE (E.G., CLUSTERING TEXT DOCUMENTS, MICROARRAY DATA)
  - Projected subspace-clustering: which dimensions to be projected on?
    - CLIQUE, ProClus
  - Feature extraction: costly and may not be effective?
  - Using frequent patterns as “features”
    - “Frequent” are inherent features
    - Mining freq. patterns may not be so expensive
- TYPICAL METHODS
  - Frequent-term-based document clustering
  - Clustering by pattern similarity in micro-array data (pClustering)

# Cluster Validity

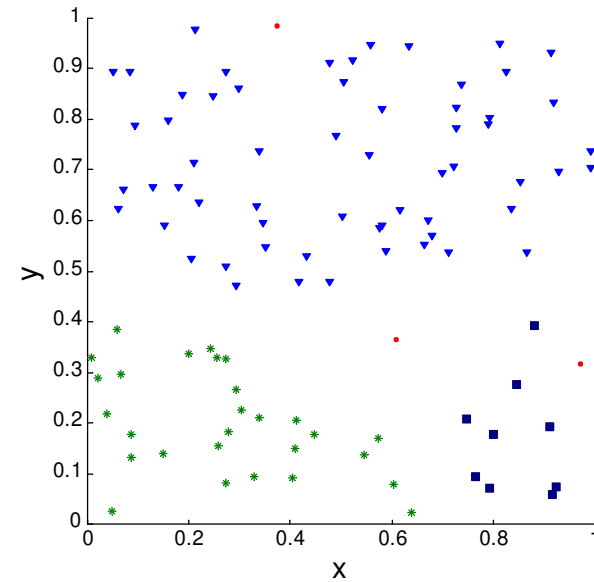
- ✗ FOR SUPERVISED CLASSIFICATION WE HAVE A VARIETY OF MEASURES TO EVALUATE HOW GOOD OUR MODEL IS
  - ✗ Accuracy, precision, recall
- ✗ FOR CLUSTER ANALYSIS, THE ANALOGOUS QUESTION IS HOW TO EVALUATE THE “GOODNESS” OF THE RESULTING CLUSTERS?
- ✗ BUT “CLUSTERS ARE IN THE EYE OF THE BEHOLDER”!
- ✗ THEN 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

# Clusters found in Random Data

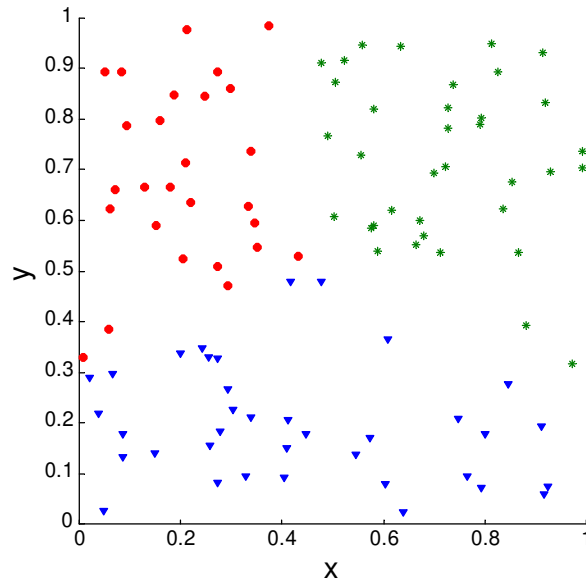
Random Points



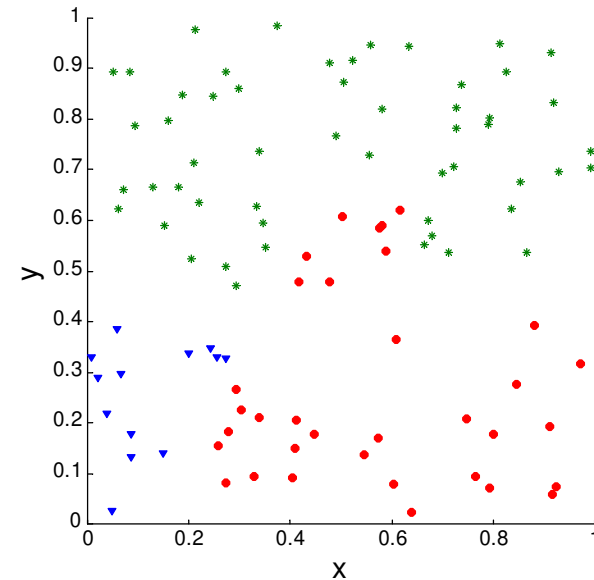
DBSCAN



K-means



Complete Link





# Different Aspects of Cluster Validation

1. DETERMINING THE **CLUSTERING TENDENCY** OF A SET OF DATA, I.E., DISTINGUISHING WHETHER NON-RANDOM STRUCTURE ACTUALLY EXISTS IN THE DATA.
2. COMPARING THE RESULTS OF A CLUSTER ANALYSIS TO EXTERNALLY KNOWN RESULTS, E.G., TO EXTERNALLY GIVEN CLASS LABELS.
3. EVALUATING HOW WELL THE RESULTS OF A CLUSTER ANALYSIS FIT THE DATA *WITHOUT* REFERENCE TO EXTERNAL INFORMATION.
  - Use only the data
4. COMPARING THE RESULTS OF TWO DIFFERENT SETS OF CLUSTER ANALYSES TO DETERMINE WHICH IS BETTER.
5. DETERMINING THE 'CORRECT' NUMBER OF CLUSTERS.

FOR 2, 3, AND 4, WE CAN FURTHER DISTINGUISH WHETHER WE WANT TO EVALUATE THE ENTIRE CLUSTERING OR JUST INDIVIDUAL CLUSTERS.

# Measures of Cluster Validity

✕ NUMERICAL MEASURES THAT ARE APPLIED TO JUDGE VARIOUS ASPECTS OF CLUSTER VALIDITY, ARE CLASSIFIED INTO THE FOLLOWING THREE TYPES.

- ✕ External Index: Used to measure the extent to which cluster labels match externally supplied class labels.

  - ✕ Entropy

- ✕ Internal Index: Used to measure the goodness of a clustering structure without respect to external information.

  - ✕ Sum of Squared Error (SSE)

- ✕ Relative Index: Used to compare two different clusterings or clusters.

  - ✕ Often an external or internal index is used for this function, e.g., SSE or entropy

✕ SOMETIMES THESE ARE REFERRED TO AS CRITERIA INSTEAD OF INDICES

- ✕ However, as a rule criterion is the general strategy and index is the numerical measure that implements the criterion.

# Measuring Cluster Validity Via Correlation

## xTWO MATRICES

### xProximity Matrix

- xDistance between any pair of rows

### xIncidence Matrix

- xOne row and one column for each data point
- xAn entry is 1 if the associated pair of points belong to the same cluster
- xAn entry is 0 if the associated pair of points belongs to different clusters

## xCOMPUTE THE CORRELATION BETWEEN THE TWO MATRICES

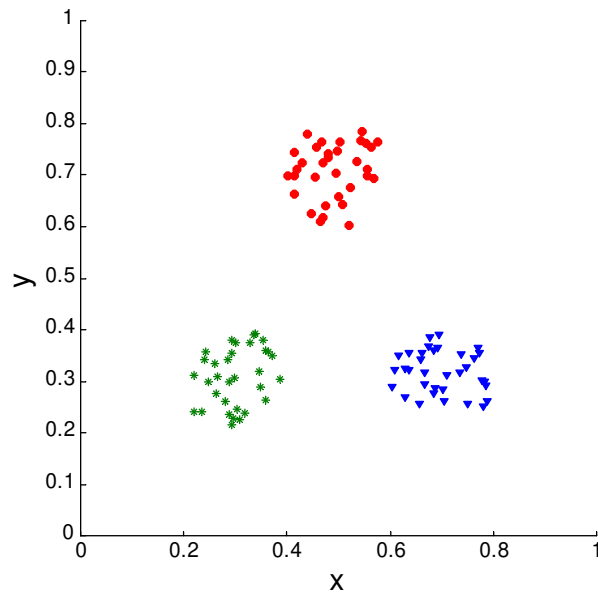
- xSince the matrices are symmetric, only the correlation between  $n(n-1)/2$  entries needs to be calculated.

## xHIGH CORRELATION INDICATES THAT POINTS THAT BELONG TO THE SAME CLUSTER ARE CLOSE TO EACH OTHER.

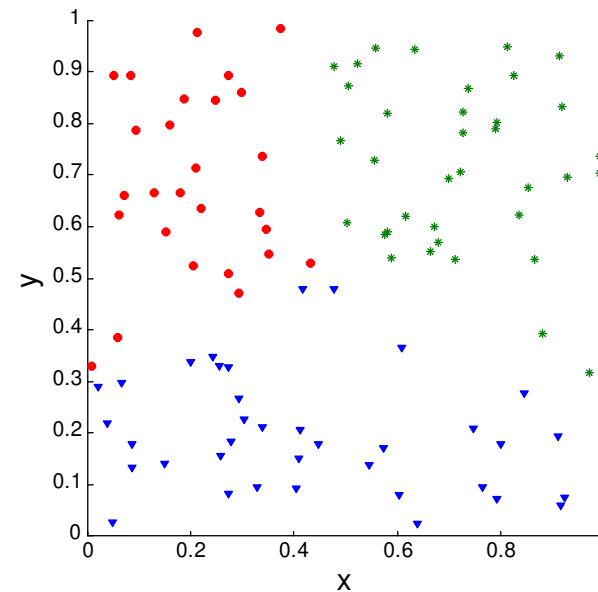
## xNOT A GOOD MEASURE FOR SOME DENSITY OR CONTIGUITY BASED CLUSTERS.

# Measuring Cluster Validity Via Correlation

✗ CORRELATION OF INCIDENCE AND PROXIMITY MATRICES FOR THE K-MEANS CLUSTERINGS OF THE FOLLOWING TWO DATA SETS.



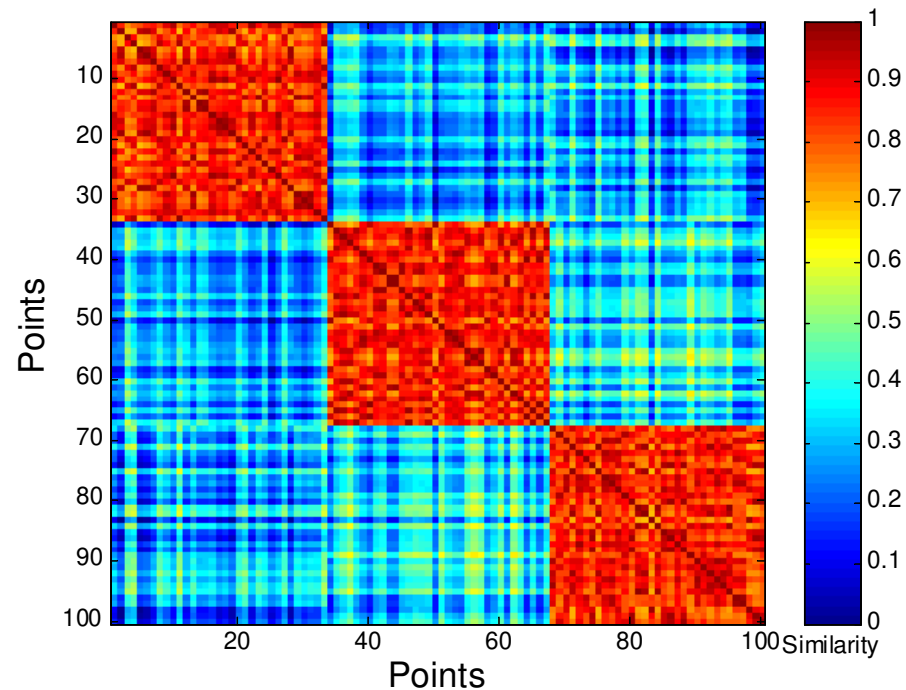
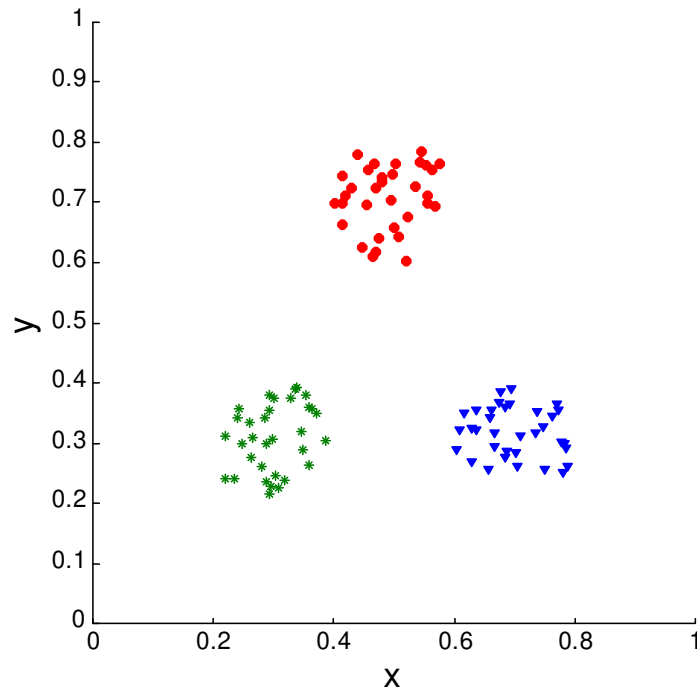
CORR = -0.9235



CORR = -0.5810

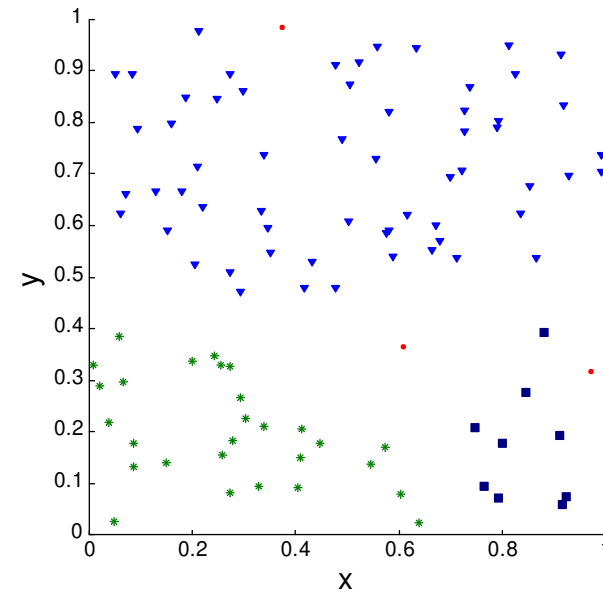
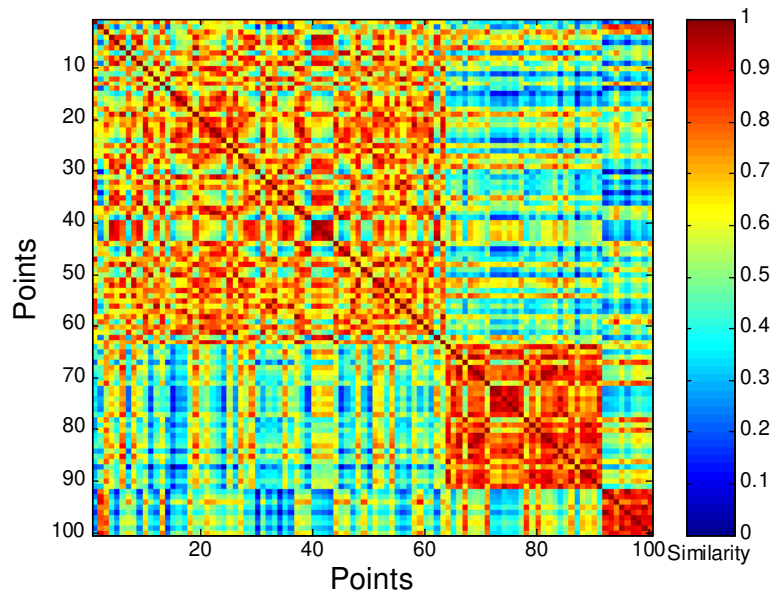
# Using Similarity Matrix for Cluster Validation

ORDER THE SIMILARITY MATRIX WITH RESPECT TO CLUSTER LABELS AND INSPECT VISUALLY.



# Using Similarity Matrix for Cluster Validation

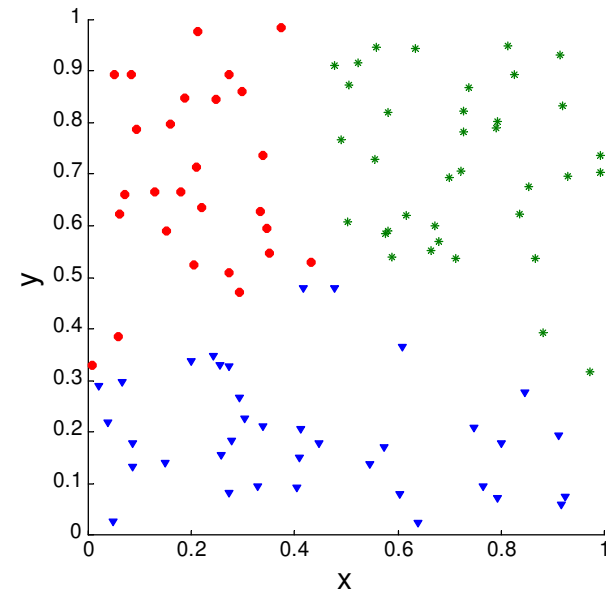
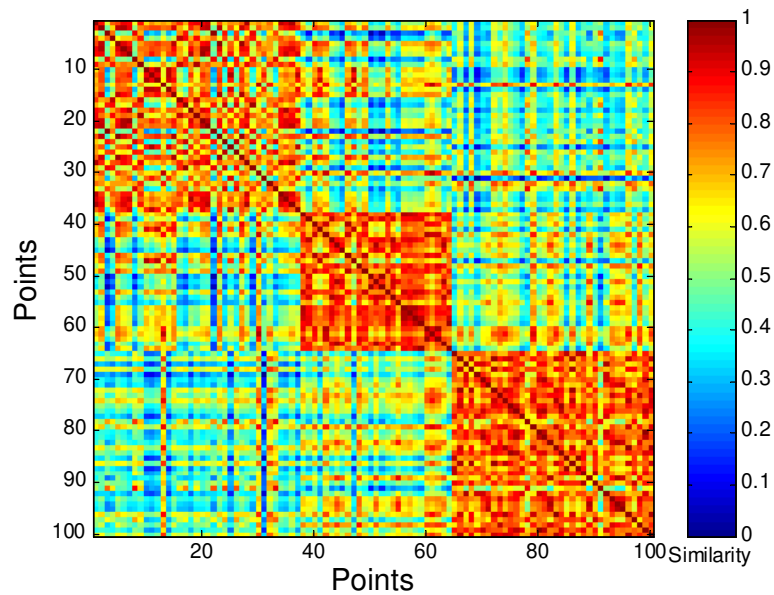
✗ CLUSTERS IN RANDOM DATA ARE NOT SO CRISP



DBSCAN

# Using Similarity Matrix for Cluster Validation

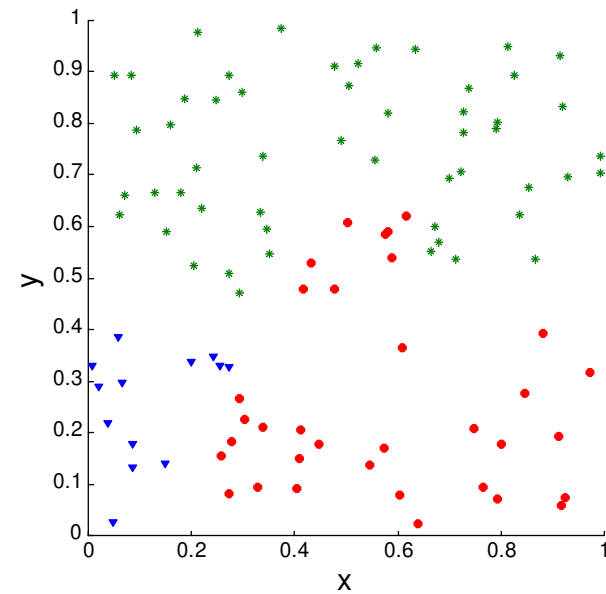
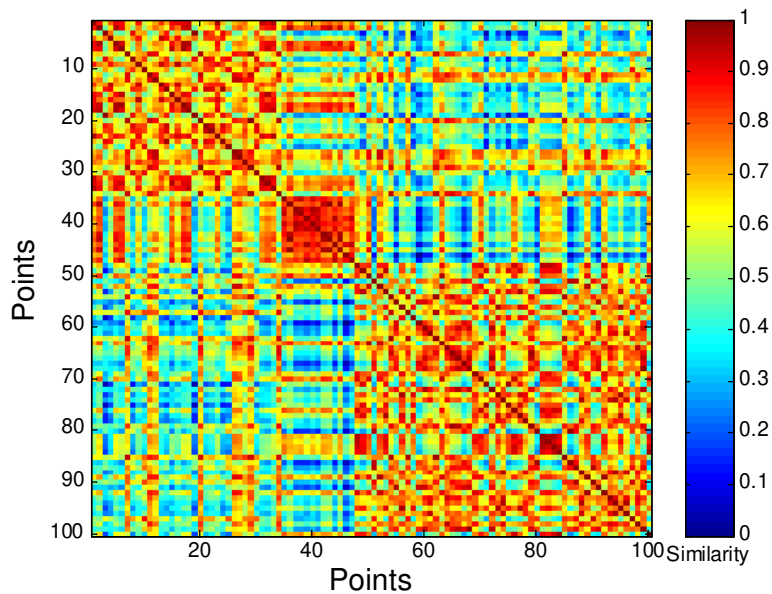
## CLUSTERS IN RANDOM DATA ARE NOT SO CRISP



## K-MEANS

# Using Similarity Matrix for Cluster Validation

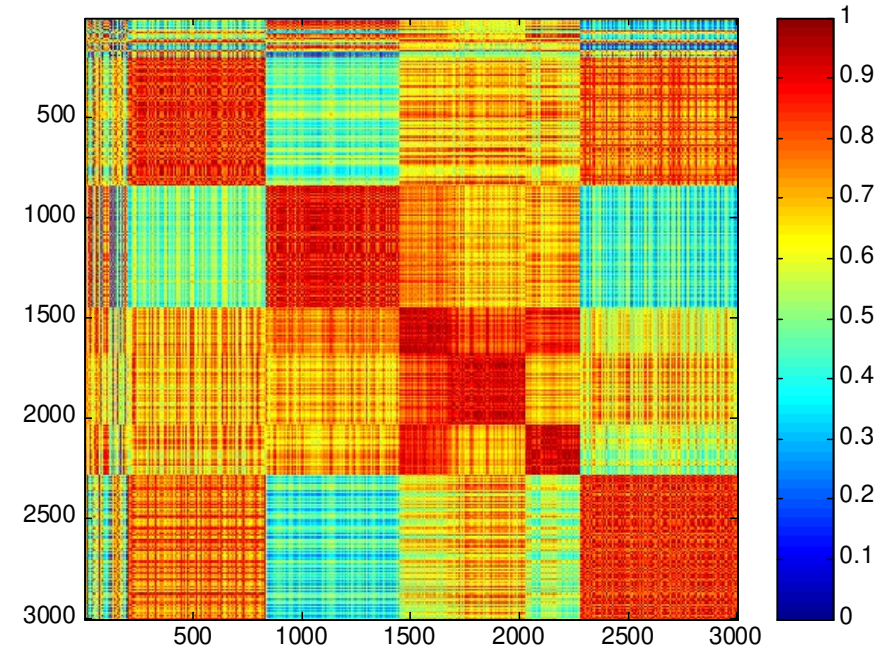
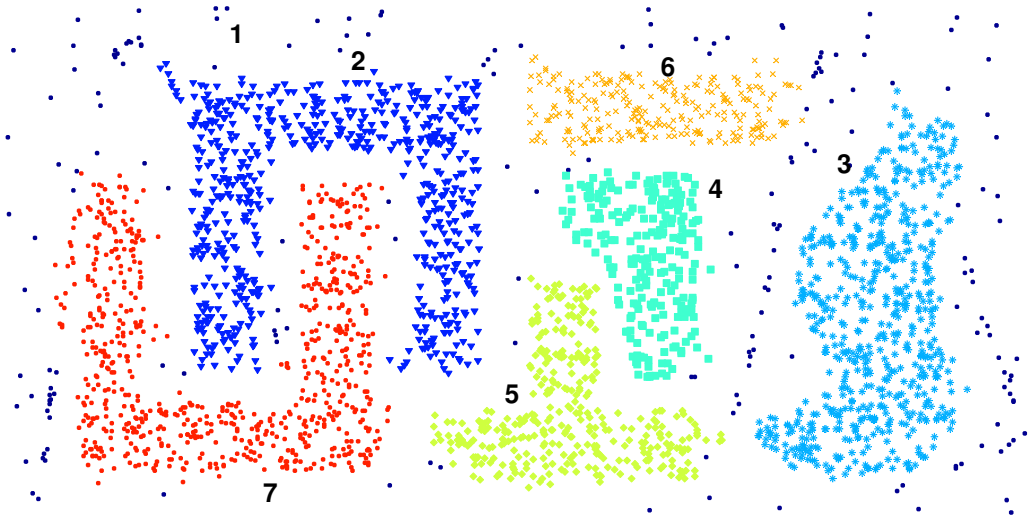
✗ CLUSTERS IN RANDOM DATA ARE NOT SO CRISP



COMPLETE LINK



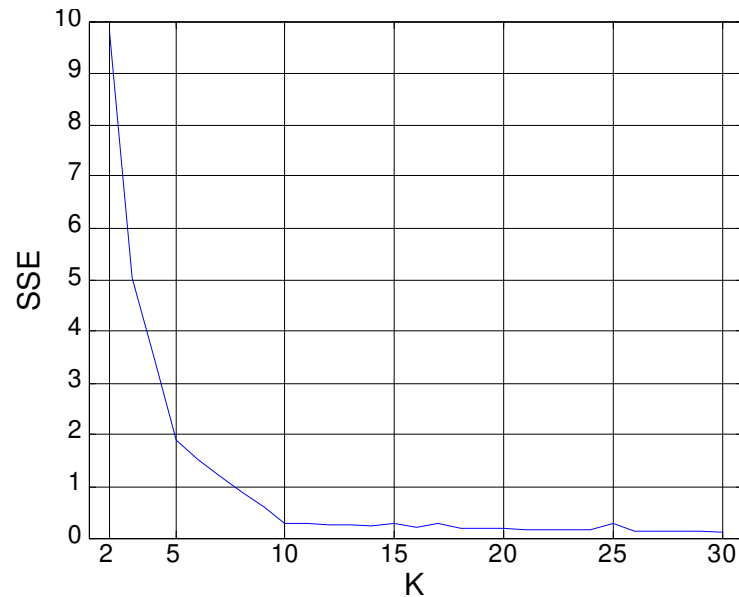
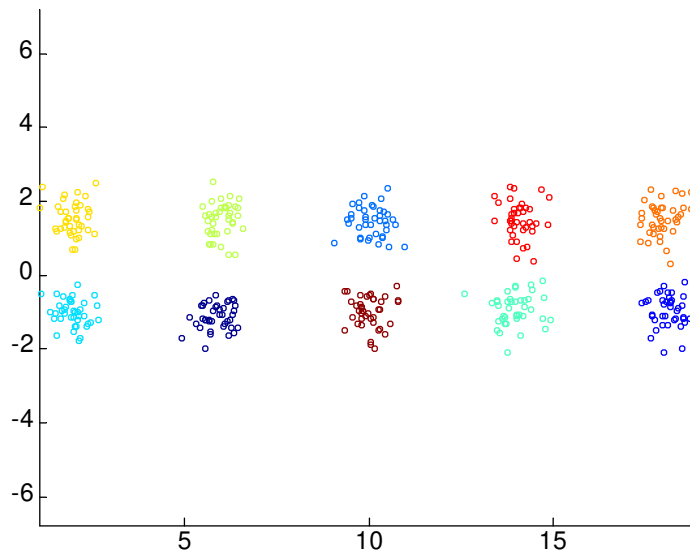
# Using Similarity Matrix for Cluster Validation



DBSCAN

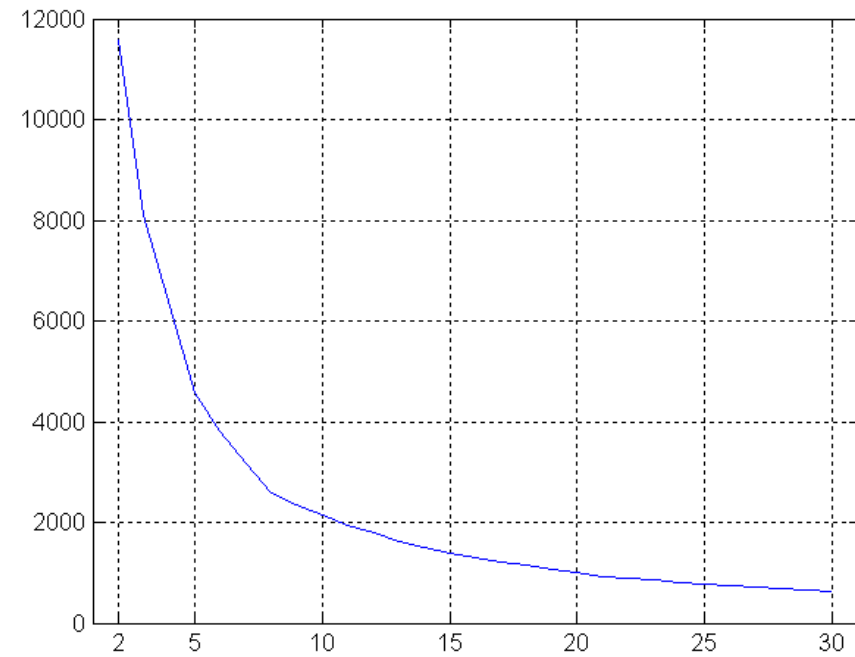
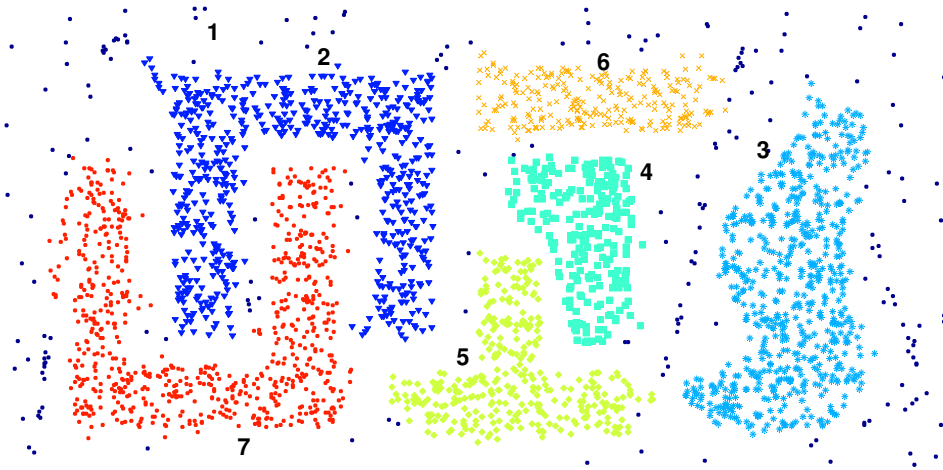
# Internal Measures: SSE

- ✗ CLUSTERS IN COMPLICATED FIGURES AREN'T WELL SEPARATED
- ✗ INTERNAL INDEX:
  - ✗ Used to measure the goodness of a clustering structure without respect to external information
- ✗ SSE IS GOOD FOR COMPARING TWO CLUSTERINGS OR TWO CLUSTERS (AVERAGE SSE).
- ✗ CAN ALSO BE USED TO ESTIMATE THE NUMBER OF CLUSTERS



# Internal Measures: SSE

✕SSE CURVE FOR A MORE COMPLICATED DATA SET



SSE OF CLUSTERS FOUND USING K-MEANS

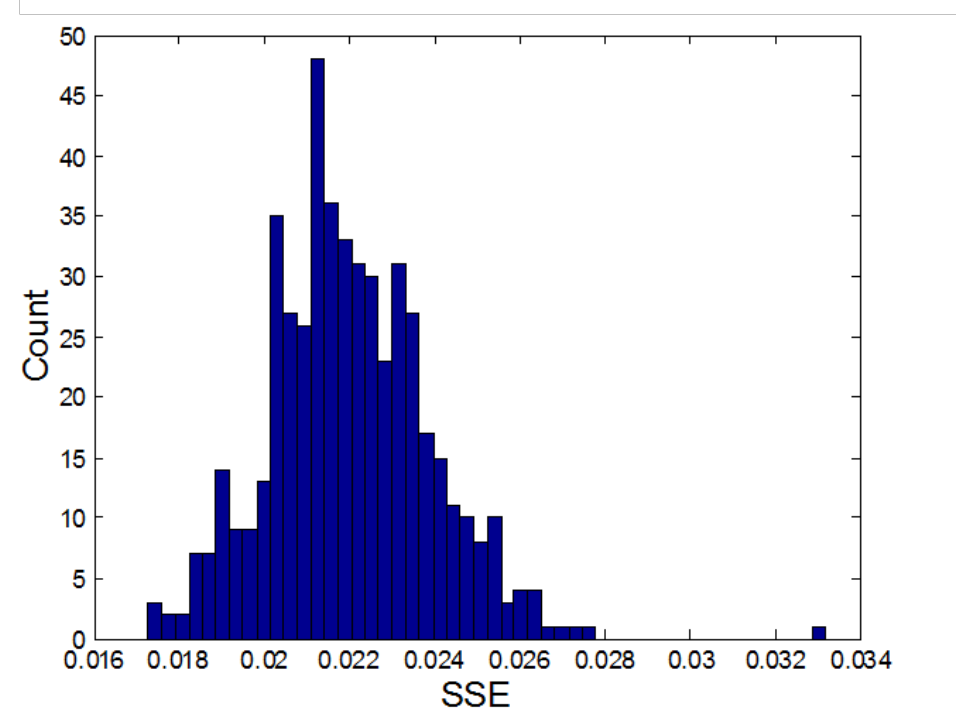
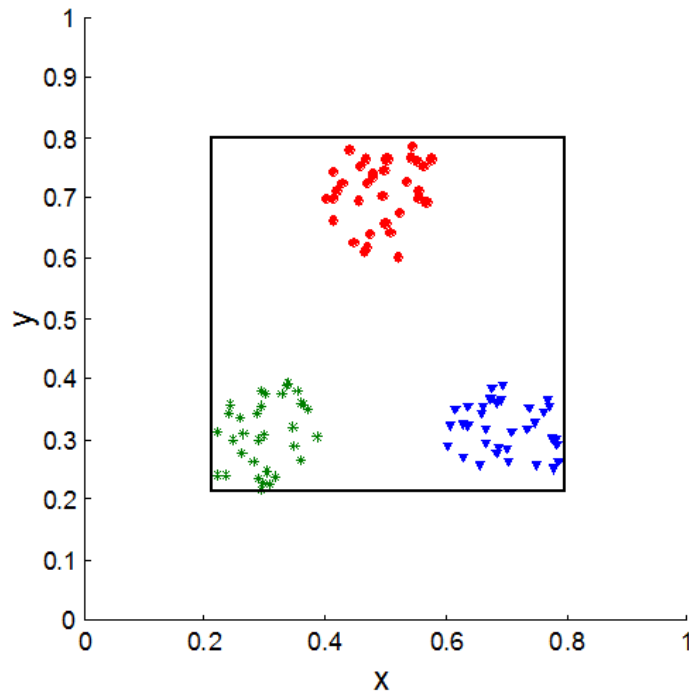
# Framework for Cluster Validity

- ✗ NEED A FRAMEWORK TO INTERPRET ANY MEASURE.
  - ✗ For example, if our measure of evaluation has the value, 10, is that good, fair, or poor?
- ✗ STATISTICS PROVIDE A FRAMEWORK FOR CLUSTER VALIDITY
  - ✗ The more “atypical” a clustering result is, the more likely it represents valid structure in the data
  - ✗ Can compare the values of an index that result from random data or clusterings to those of a clustering result.
    - ✗ If the value of the index is unlikely, then the cluster results are valid
  - ✗ These approaches are more complicated and harder to understand.
- ✗ FOR COMPARING THE RESULTS OF TWO DIFFERENT SETS OF CLUSTER ANALYSES, A FRAMEWORK IS LESS NECESSARY.
  - ✗ However, there is the question of whether the difference between two index values is significant

# Statistical Framework for SSE

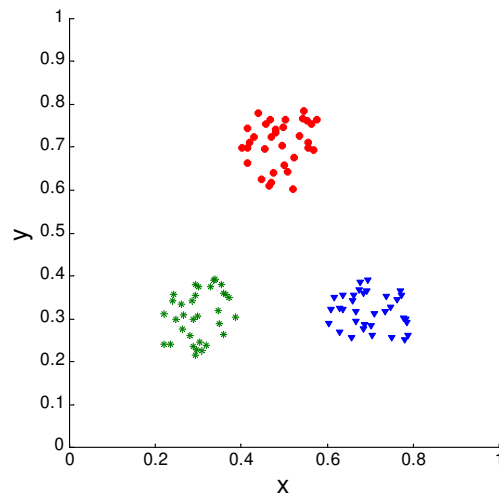
## ✕EXAMPLE

- ✕Compare SSE of 0.005 against three clusters in random data
- ✕Histogram shows SSE of three clusters in 500 sets of random data points of size 100 distributed over the range 0.2 – 0.8 for x and y values

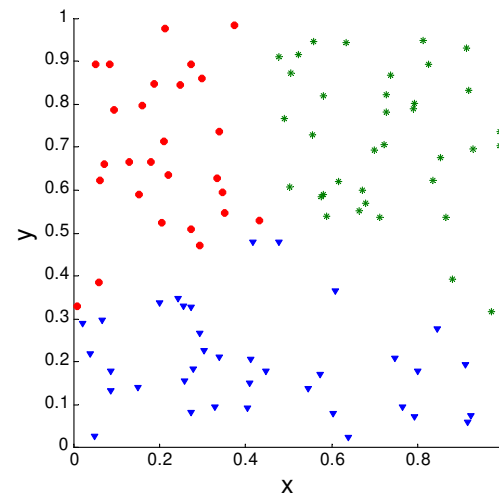


# Statistical Framework for Correlation

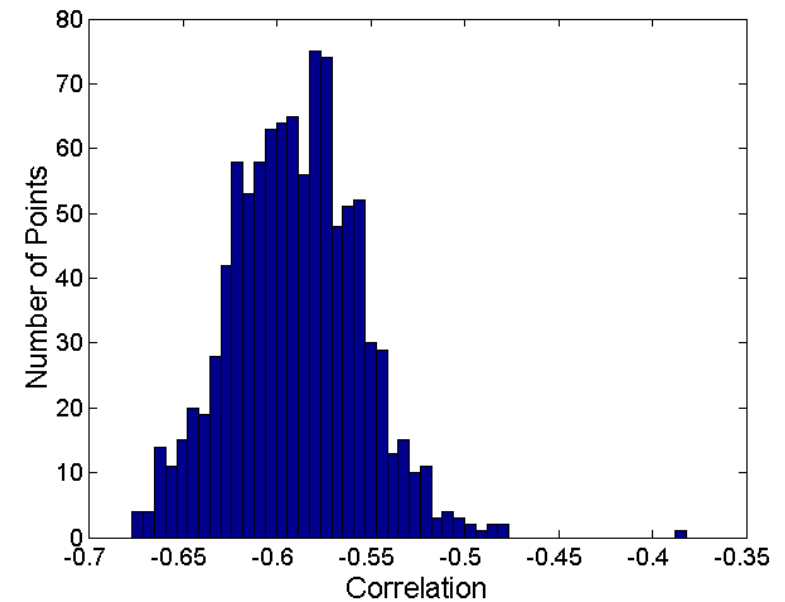
✕CORRELATION OF INCIDENCE AND PROXIMITY MATRICES FOR THE K-MEANS CLUSTERINGS OF THE FOLLOWING TWO DATA SETS.



CORR = -0.9235



CORR = -0.5810



# 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

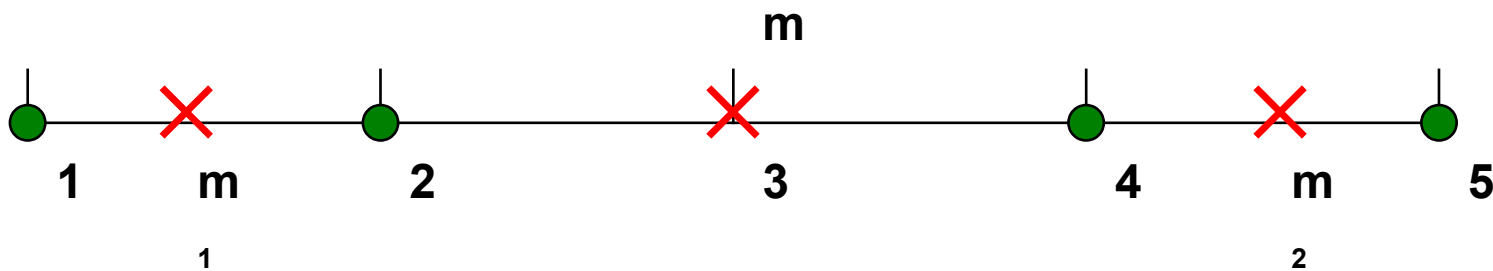
$$BSS = \sum_i |C_i| (m - m_i)^2$$

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

# Internal Measures: Cohesion and Separation

✗EXAMPLE: SSE

✗BSS + WSS = constant



K=1 CLUSTER:

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

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

$$Total = 10 + 0 = 10$$

K=2 CLUSTERS:

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

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

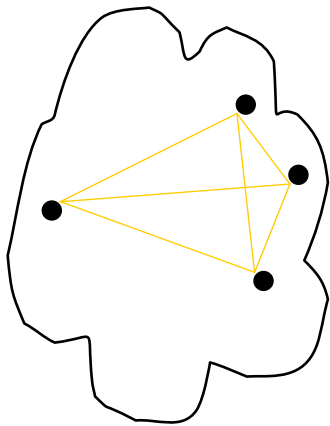
$$Total = 1 + 9 = 10$$



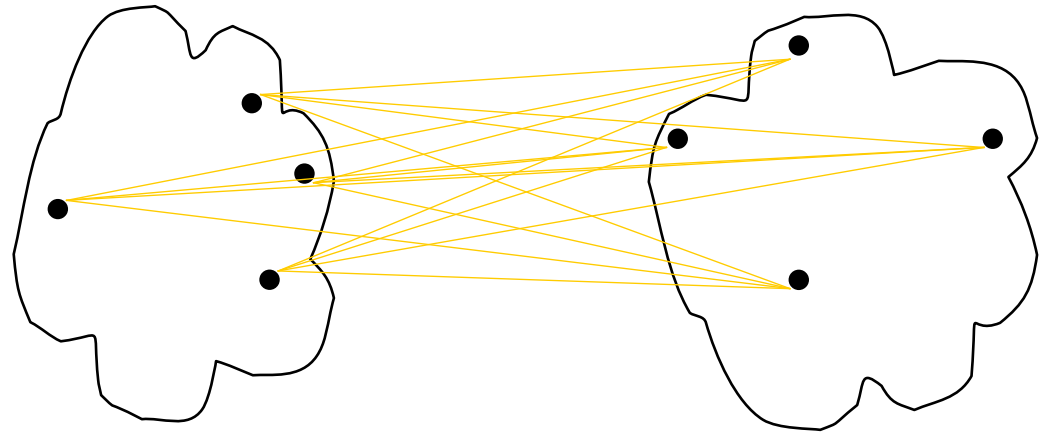
# Internal Measures: Cohesion and Separation

✕ A PROXIMITY GRAPH BASED APPROACH CAN ALSO BE USED FOR COHESION AND SEPARATION.

- ✕ Cluster cohesion is the sum of the weight of links within a cluster.
- ✕ Cluster separation is the sum of the weights between nodes in the cluster and nodes outside the cluster.



COHESION



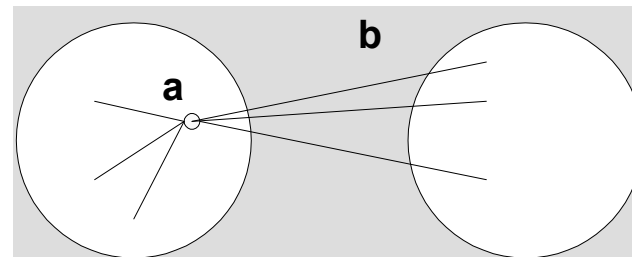
SEPARATION

# Internal Measures: Silhouette Coefficient

- ✗ SILHOUETTE COEFFICIENT COMBINE IDEAS OF BOTH COHESION AND SEPARATION, BUT FOR INDIVIDUAL POINTS, AS WELL AS CLUSTERS AND CLUSTERINGS
- ✗ FOR AN INDIVIDUAL POINT, I
  - ✗ Calculate  $a$  = average distance of  $i$  to the points in its cluster
  - ✗ Calculate  $b$  = min (average distance of  $i$  to points in another cluster)
  - ✗ The silhouette coefficient for a point is then given by

$$s = 1 - a/b \quad \text{if } a < b, \quad (\text{or } s = b/a - 1 \text{ if } a \geq b, \text{ not the usual case})$$

- ✗ Typically between 0 and 1.
- ✗ The closer to 1 the better.



- ✗ CAN CALCULATE THE AVERAGE SILHOUETTE WIDTH FOR A CLUSTER OR A CLUSTERING

# External Measures of Cluster Validity: Entropy and Purity

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

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

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

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

# More measures

Table I  
INTERNAL CLUSTERING VALIDATION MEASURES

Measure	Notation	Definition	Optimal value
1 Root-mean-square std dev	$RMSSTD$	$\{\sum_i \sum_{x \in C_i} \ x - c_i\ ^2 / [P \sum_i (n_i - 1)]\}^{\frac{1}{2}}$	Elbow
2 R-squared	$RS$	$(\sum_{x \in D} \ x - c\ ^2 - \sum_i \sum_{x \in C_i} \ x - c_i\ ^2) / \sum_{x \in D} \ x - c\ ^2$	Elbow
3 Modified Hubert $\Gamma$ statistic	$\Gamma$	$\frac{2}{n(n-1)} \sum_{x \in D} \sum_{y \in D} d(x, y) d_{x \in C_i, y \in C_j}(c_i, c_j)$	Elbow
4 Calinski-Harabasz index	$CH$	$\frac{\sum_i n_i d^2(c_i, c) / (NC - 1)}{\sum_i \sum_{x \in C_i} d^2(x, c_i) / (n - NC)}$	Max
5 $I$ index	$I$	$(\frac{1}{NC} \cdot \frac{\sum_{x \in D} d(x, c)}{\sum_i \sum_{x \in C_i} d(x, c_i)} \cdot \max_{i,j} d(c_i, c_j))^p$	Max
6 Dunn's indices	$D$	$\min_i \{ \min_j ( \frac{\min_{x \in C_i, y \in C_j} d(x, y)}{\max_k \{ \max_{x, y \in C_k} d(x, y) \}} ) \}$	Max
7 Silhouette index	$S$	$\frac{1}{NC} \sum_i \{ \frac{1}{n_i} \sum_{x \in C_i} \frac{b(x) - a(x)}{\max[b(x), a(x)]} \}$ $a(x) = \frac{1}{n_i - 1} \sum_{y \in C_i, y \neq x} d(x, y), b(x) = \min_{j, j \neq i} [ \frac{1}{n_j} \sum_{y \in C_j} d(x, y) ]$	Max
8 Davies-Bouldin index	$DB$	$\frac{1}{NC} \sum_i \max_{j, j \neq i} \{ [ \frac{1}{n_i} \sum_{x \in C_i} d(x, c_i) + \frac{1}{n_j} \sum_{x \in C_j} d(x, c_j) ] / d(c_i, c_j) \}$	Min
9 Xie-Beni index	$XB$	$[ \sum_i \sum_{x \in C_i} d^2(x, c_i) ] / [ n \cdot \min_{i, j \neq i} d^2(c_i, c_j) ]$	Min
10 SD validity index	$SD$	$Dis(NC_{max}) Scat(NC) + Dis(NC)$ $Scat(NC) = \frac{1}{NC} \sum_i \  \sigma(C_i) \  / \  \sigma(D) \ , Dis(NC) = \frac{\max_{i,j} d(c_i, c_j)}{\min_{i,j} d(c_i, c_j)} \sum_i (\sum_j d(c_i, c_j))^{-1}$	Min
11 S_Dbw validity index	$S\_Dbw$	$Scat(NC) + Dens\_bw(NC)$ $Dens\_bw(NC) = \frac{1}{NC(NC-1)} \sum_i [ \sum_{j, j \neq i} \frac{\sum_{x \in C_i \cup C_j} f(x, u_{ij})}{\max \{ \sum_{x \in C_i} f(x, c_i), \sum_{x \in C_j} f(x, c_j) \}} ]$	Min

$D$ : data set;  $n$ : number of objects in  $D$ ;  $c$ : center of  $D$ ;  $P$ : attributes number of  $D$ ;  $NC$ : number of clusters;  $C_i$ : the  $i$ -th cluster;  $n_i$ : number of objects in  $C_i$ ;  $c_i$ : center of  $C_i$ ;  $\sigma(C_i)$ : variance vector of  $C_i$ ;  $d(x, y)$ : distance between  $x$  and  $y$ ;  $\|X_i\| = (X_i^T \cdot X_i)^{\frac{1}{2}}$

# 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*