# 12th Summer School on Data Science



## Gerasimos (Jerry) Spanakis

Most slides are based on:

Introduction to Data Mining (by Tan, Steinbach, Karpatne, Kumar)
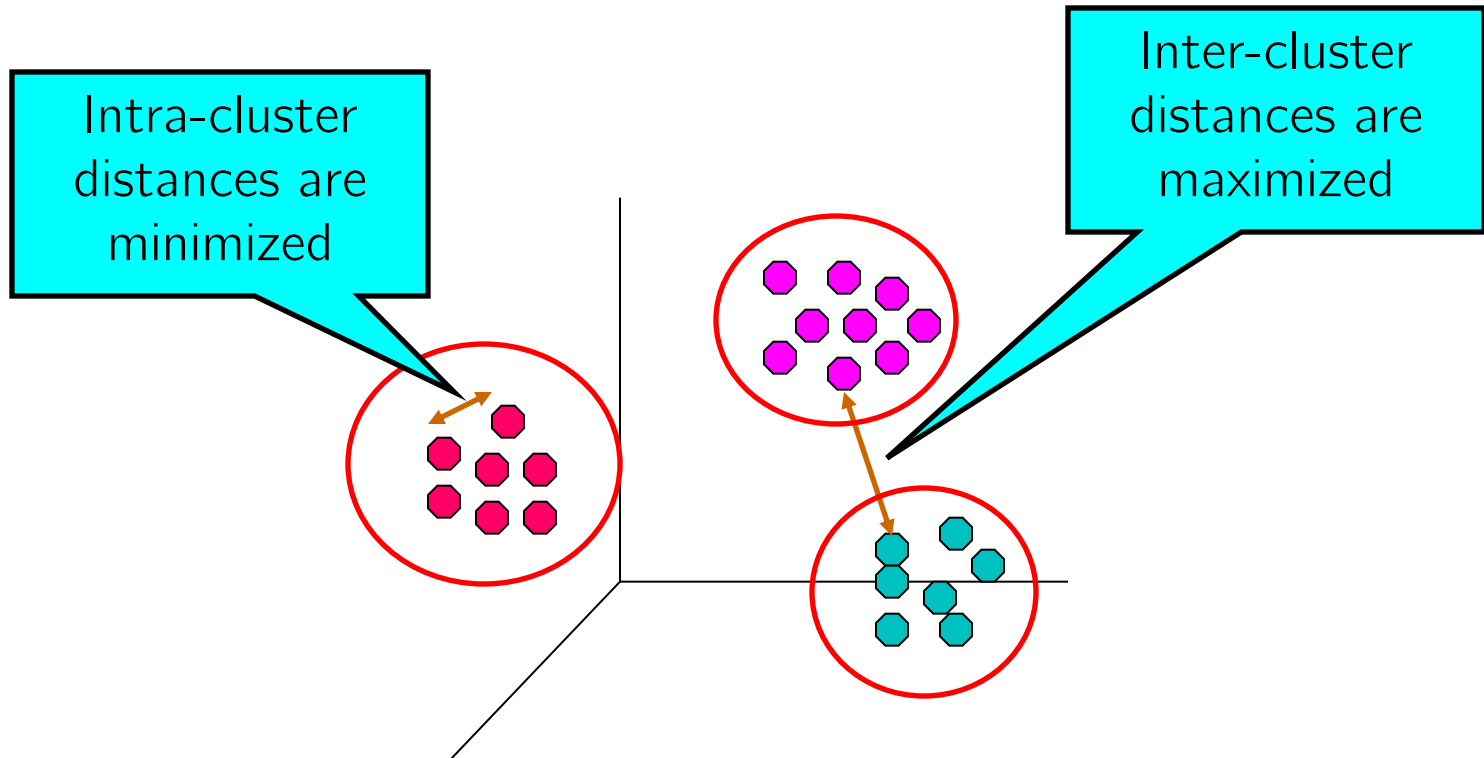
Most memes are taken from giphy.com

# Agenda

- What is clustering and why is challenging?
- Algorithms for Clustering
    - K-means
    - Hierarchical Clustering
    - DBSCAN
- Clustering Validation
    - How to make sure your clustering makes sense?
- Lab on simple clustering tasks using WEKA

# What is clustering?

- 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**

Intra-cluster distances are minimized

Inter-cluster distances are maximized

# Easy, peasy, right?



YOU on your first clustering attempt

YOU chasing a "good clustering"

NO, I'M SERIOUS

# Examples of clustering applications

- Information retrieval: Document clustering
- Marketing: Discover distinct groups in customer bases (e.g. facebook grouping: "People established adult life")
- Land use: Areas of similar land use in earth observation database
- Insurance: Groups of policy holders with a high average claim cost
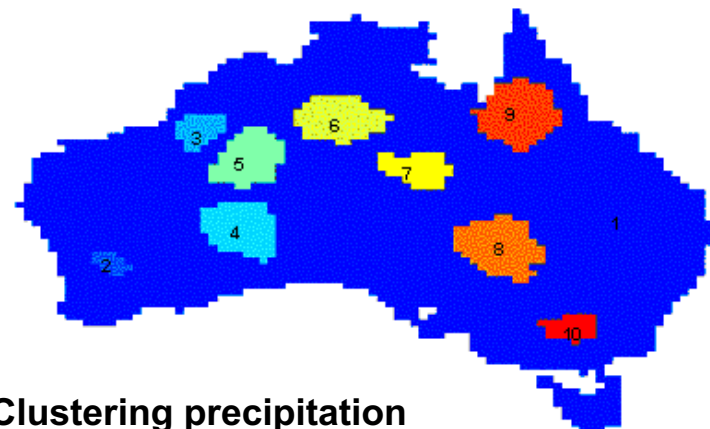- ...

# Why do it?

- ## Understanding

  – Group related documents for browsing

  – Group genes and proteins that have similar functionality

  – Group stocks with similar price fluctuations

- ## Summarization
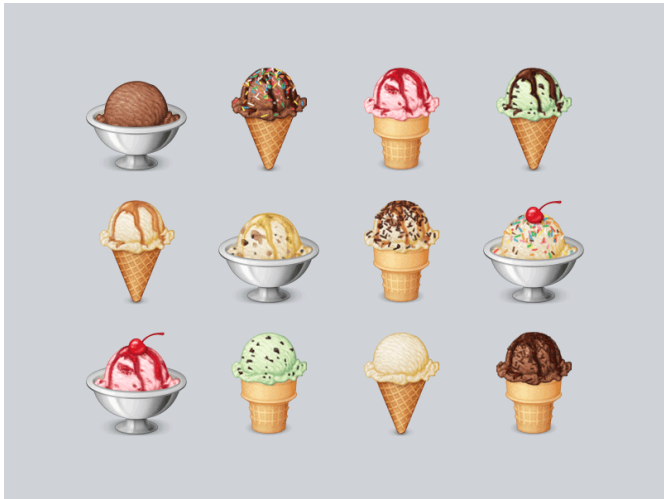
  – Reduce the size of large data sets

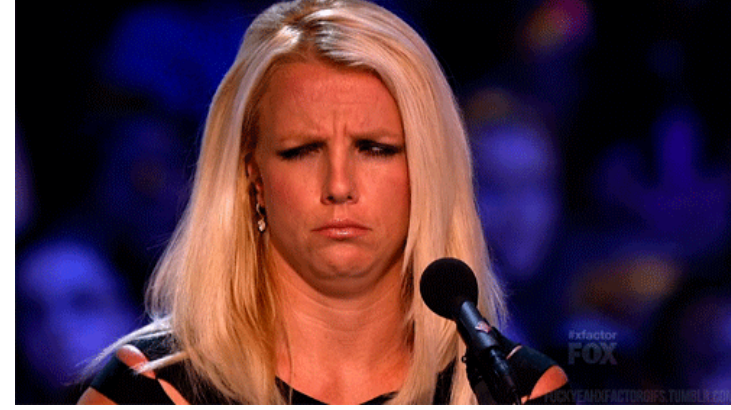| | Discovered Clusters | Industry Group |
|---|---|---|
| **1** | Applied-Matl-DOWN,Bay-Network-Down,3-COM-DOWN, Cabletron-Sys-DOWN,CISCO-DOWN,HP-DOWN, DSC-Comm-DOWN,INTEL-DOWN,LSI-Logic-DOWN, Micron-Tech-DOWN,Texas-Inst-Down,Tellabs-Inc-Down, Natl-Semiconduct-DOWN,Oracl-DOWN,SGI-DOWN, Sun-DOWN | Technology1-DOWN |
| **2** | 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 |
| **3** | Fannie-Mae-DOWN,Fed-Home-Loan-DOWN, MBNA-Corp-DOWN,Morgan-Stanley-DOWN | Financial-DOWN |
| **4** | Baker-Hughes-UP,Dresser-Inds-UP,Halliburton-HLD-UP, Louisiana-Land-UP,Phillips-Petro-UP,Unocal-UP, Schlumberger-UP | Oil-UP |

**Clustering precipitation in Australia**

# Why do it (again)

- "The revolution (in AI) will not be supervised"
- We need to have models that understand the world, like humans do
  - We don't give many "labels" to humans
  - They just learn by observing the world

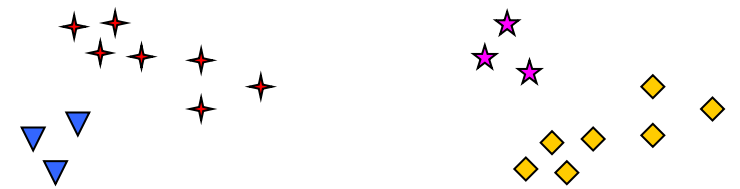# Clustering IS Ambiguous



How many clusters?

Six Clusters

Two Clusters

Four Clusters

# What is a good clustering?

- A <u>good clustering</u> method will produce high quality clusters

  – high <u>intra-class</u> similarity: <span style="color:green">cohesive</span> within clusters

  – low <u>inter-class</u> similarity: <span style="color:green">distinctive</span> between clusters

- The <u>quality</u> of a clustering method depends on

  – Data, distance, ... (see next slide)

  – its implementation,

  – Its ability to discover some or all of the <u>hidden</u> patterns

# Input data matters



- Type of data in the input
  - Measurements?
  - Image? Text? Timeseries?
- Type of distance used
  - Central to clustering
  - Depends on data and application
- Data characteristics that affect proximity and/or density are
  - Dimensionality (issues with sparseness)
  - Attribute type
  - Special relationships in the data
- Noise and Outliers
  - Often interfere with the operation of the clustering algorithm

# Major Clustering Approaches

- Partitioning approach:
  - Construct various partitions and then evaluate them by some criterion, e.g., minimizing the sum of square errors
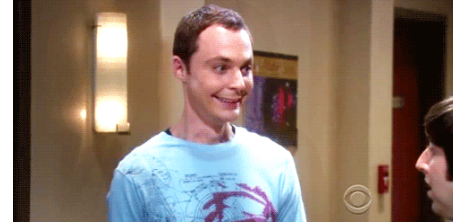  - We will see: **K-means**
- Hierarchical approach:
  - Create a hierarchical decomposition of the set of data (or objects) using some criterion
  - We will see: **Agglomerative Clustering**
- Density-based approach:
  - Based on connectivity and density functions
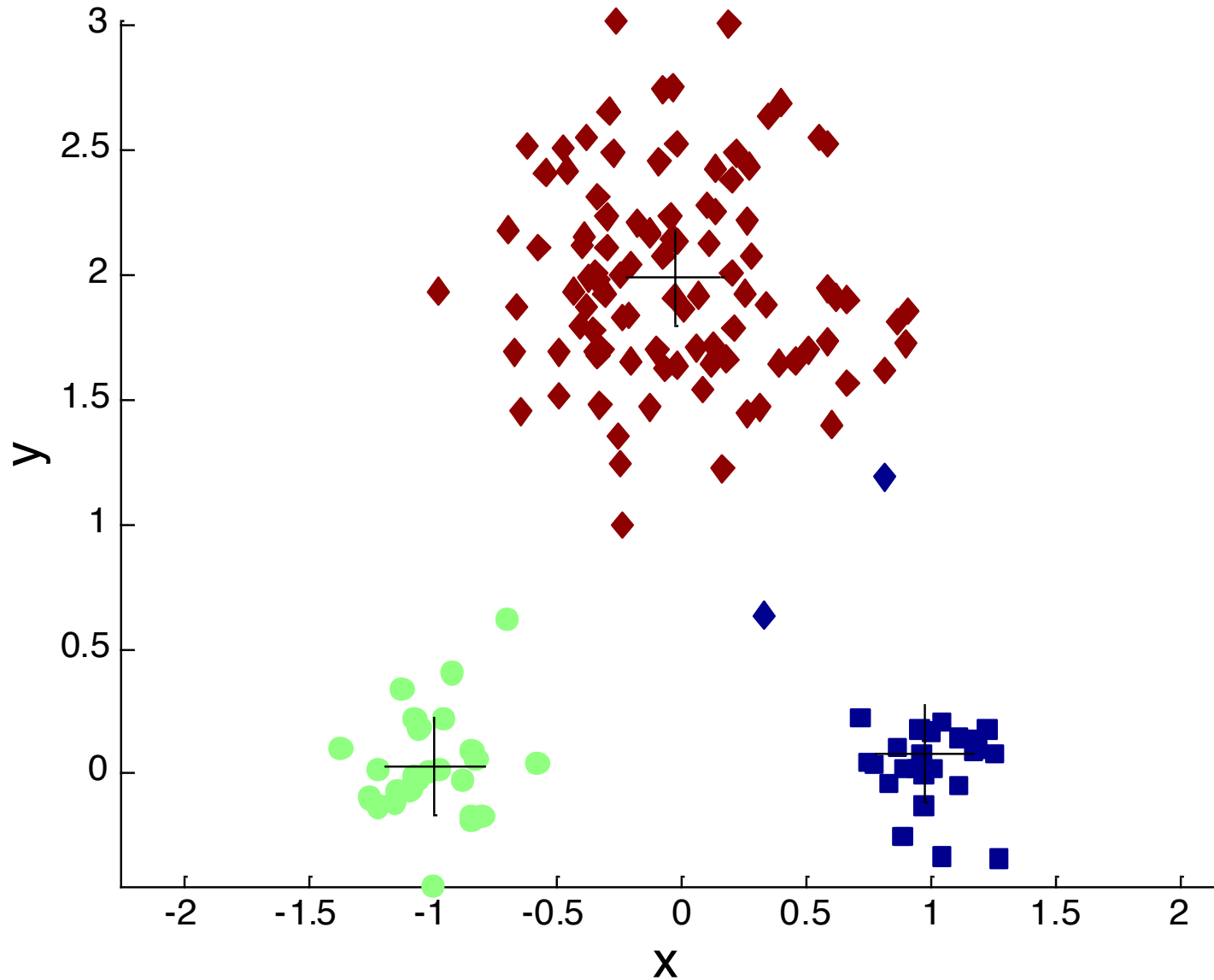  - We will see: **DBSCAN**

# K-means Clustering

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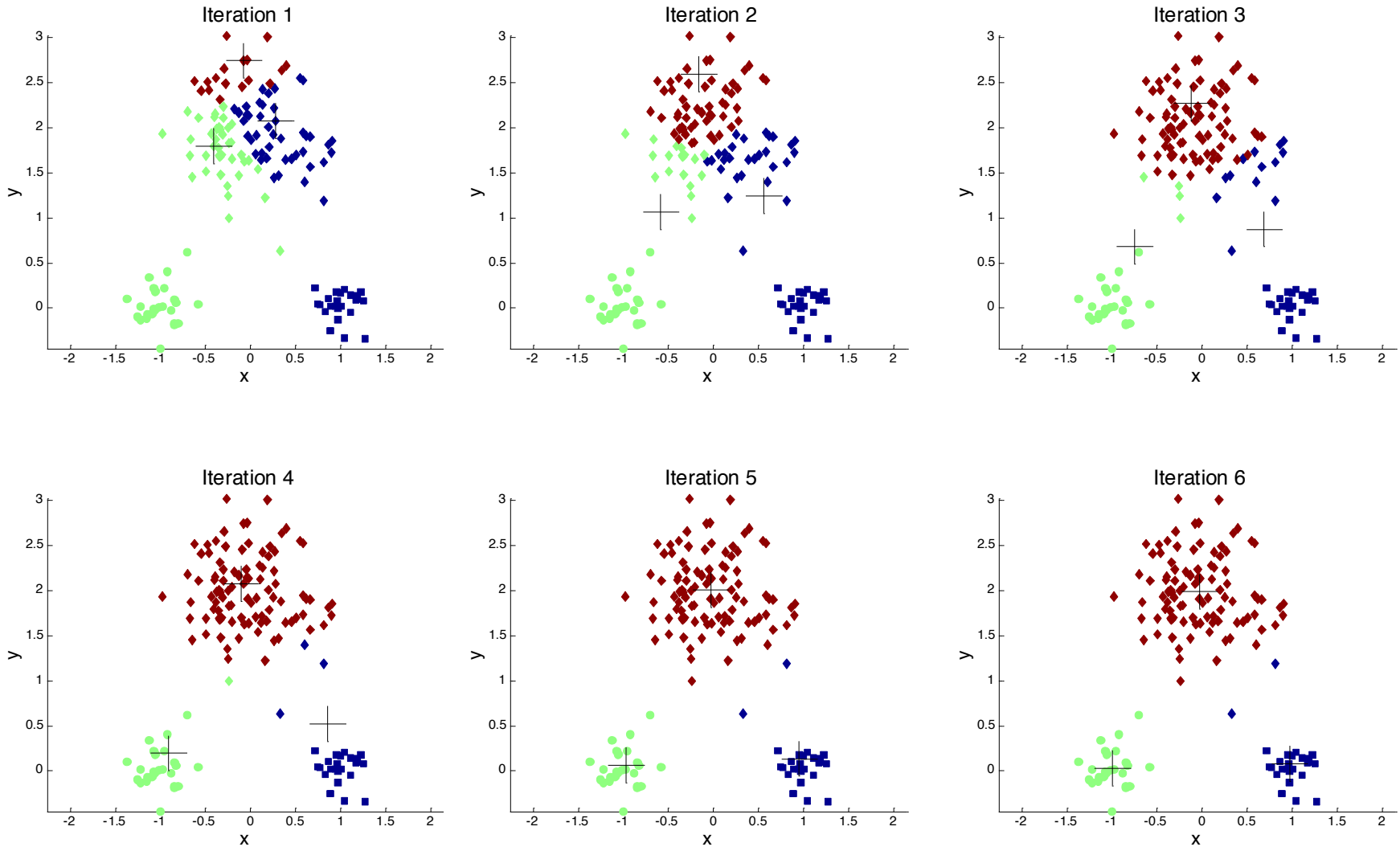- Each point is assigned to the cluster with the closest centroid

---

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

---

# Example of K-means Clustering



Iteration 6

# Example of K-means Clustering



Iteration 1

Iteration 2

Iteration 3

Iteration 4

Iteration 5

Iteration 6

# K-means Clustering - Facts

- 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.

- We need a distance measure:
  Euclidean, cosine, correlation, etc.

- K-means will converge after a few iterations

  – Often the stopping condition is changed to 'Until relatively few points change clusters'

# Evaluating K-means Clusters

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

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

  - $x$ is a data point in cluster $C_i$ and $m_i$ is the representative point for cluster $C_i$
    - can show that $m_i$ corresponds to the center (mean) of the cluster
  - Given two sets of clusters, we prefer 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
- How do we select the initial centroids?

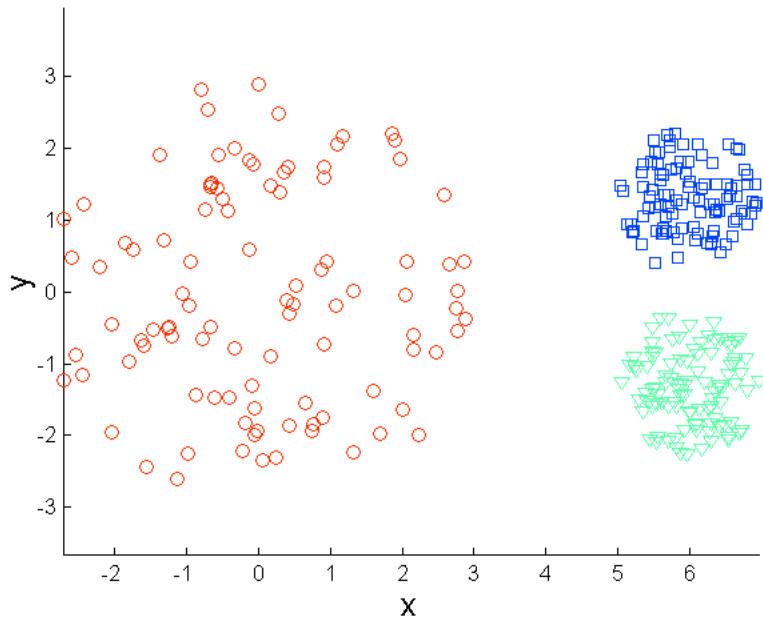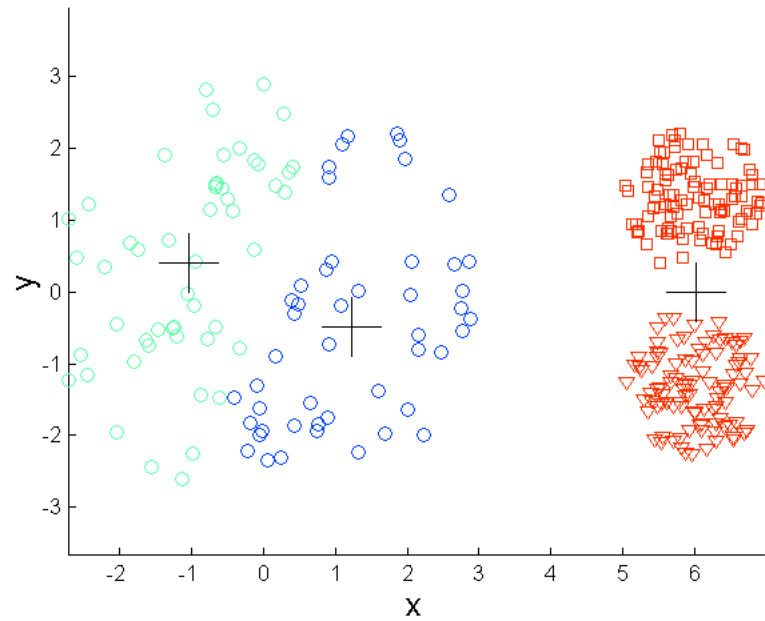# Limitations of K-means: Differing Sizes



**Original Points**

**K-means (3 Clusters)**

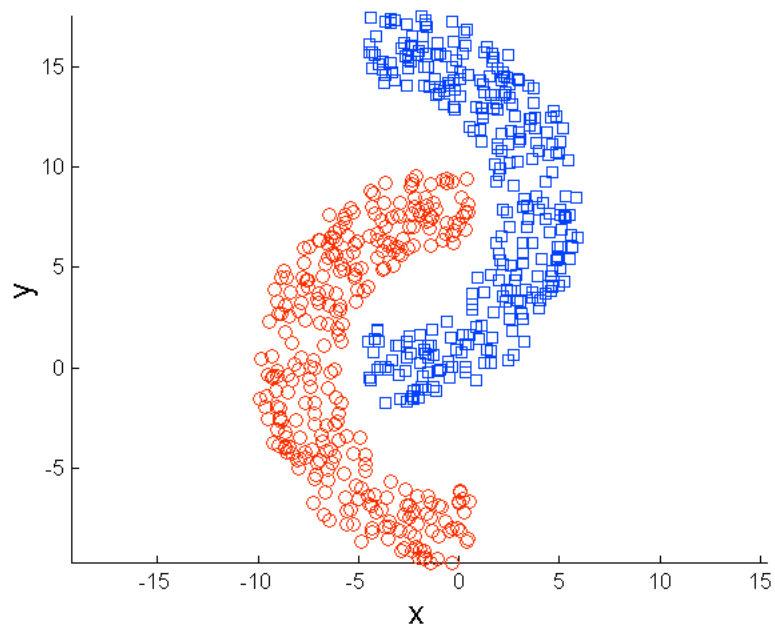# Limitations of K-means: Differing Density



**Original Points**

**K-means (3 Clusters)**
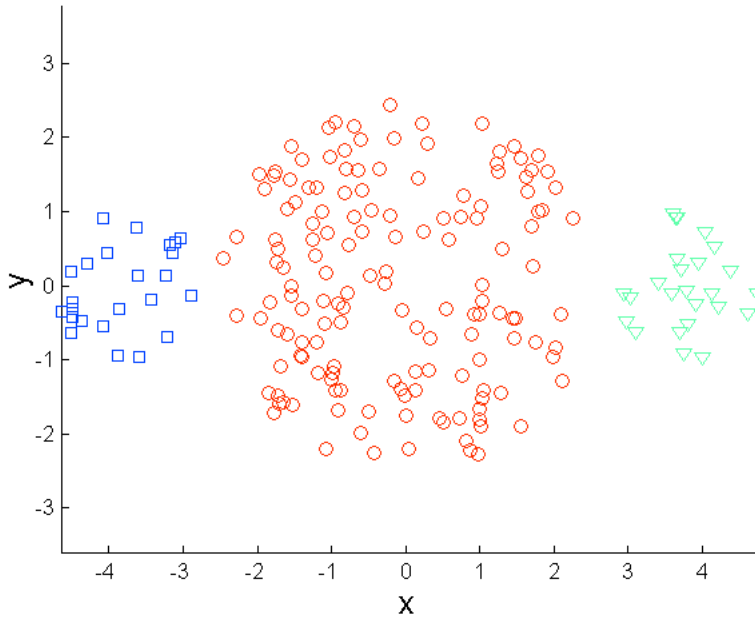
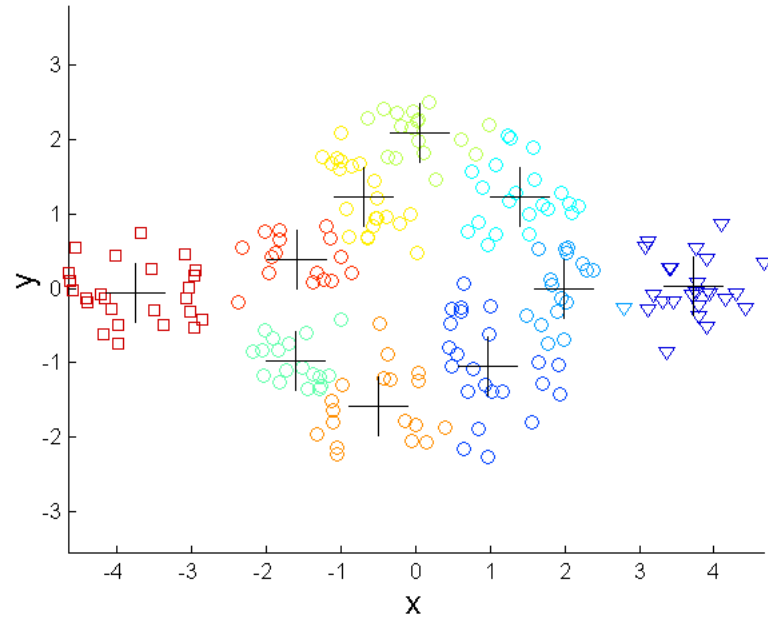# Limitations of K-means: (Non-globular) Shapes



**Original Points**

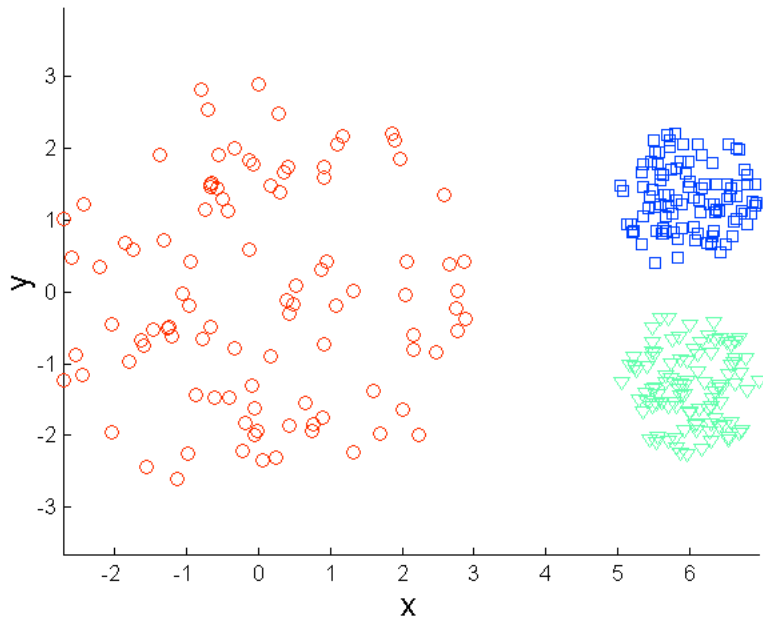**K-means (2 Clusters)**
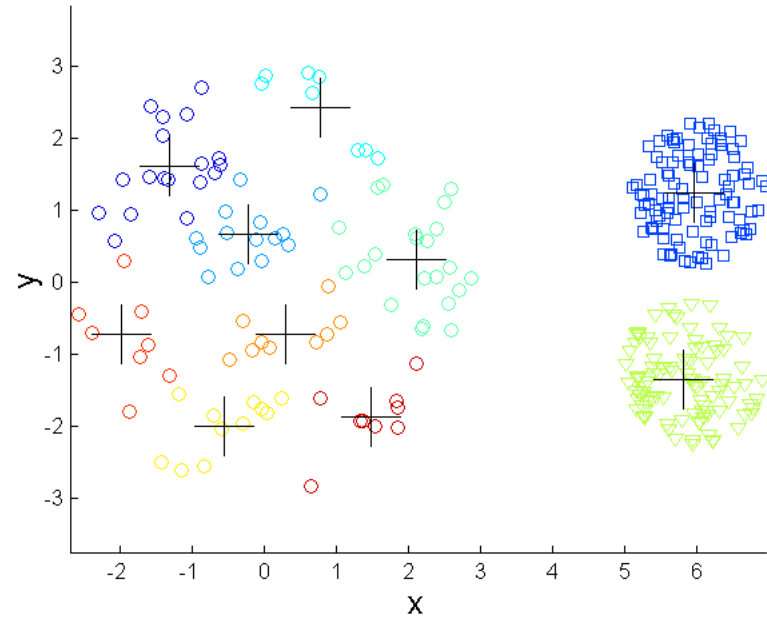
# Overcoming K-means Limitations



Original Points

K-means Clusters

One solution is to use many clusters.
    Find parts of clusters, but need to put together.
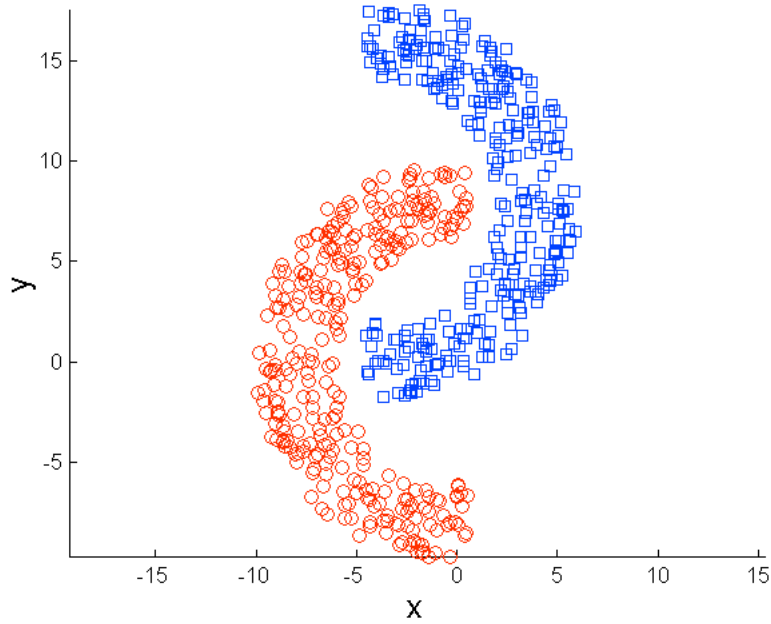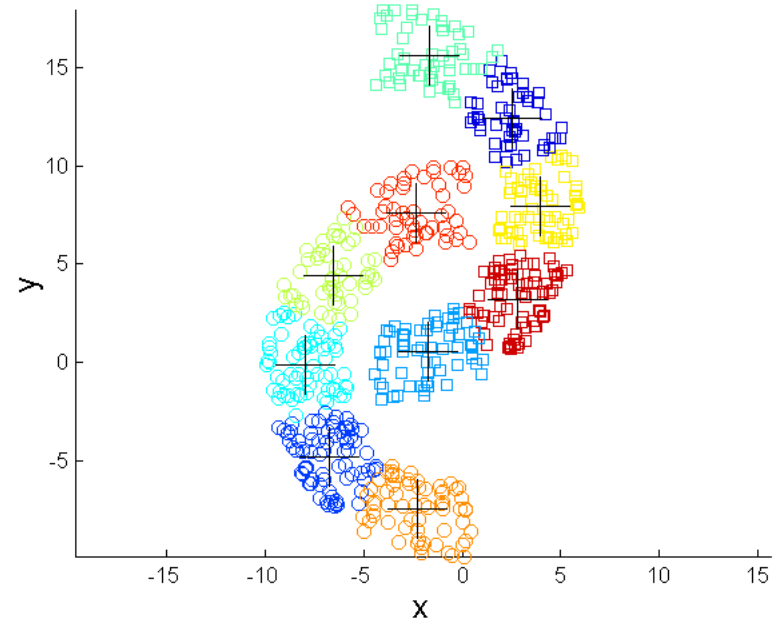
# Overcoming K-means Limitations



**Original Points**

**K-means Clusters**

# Overcoming K-means Limitations
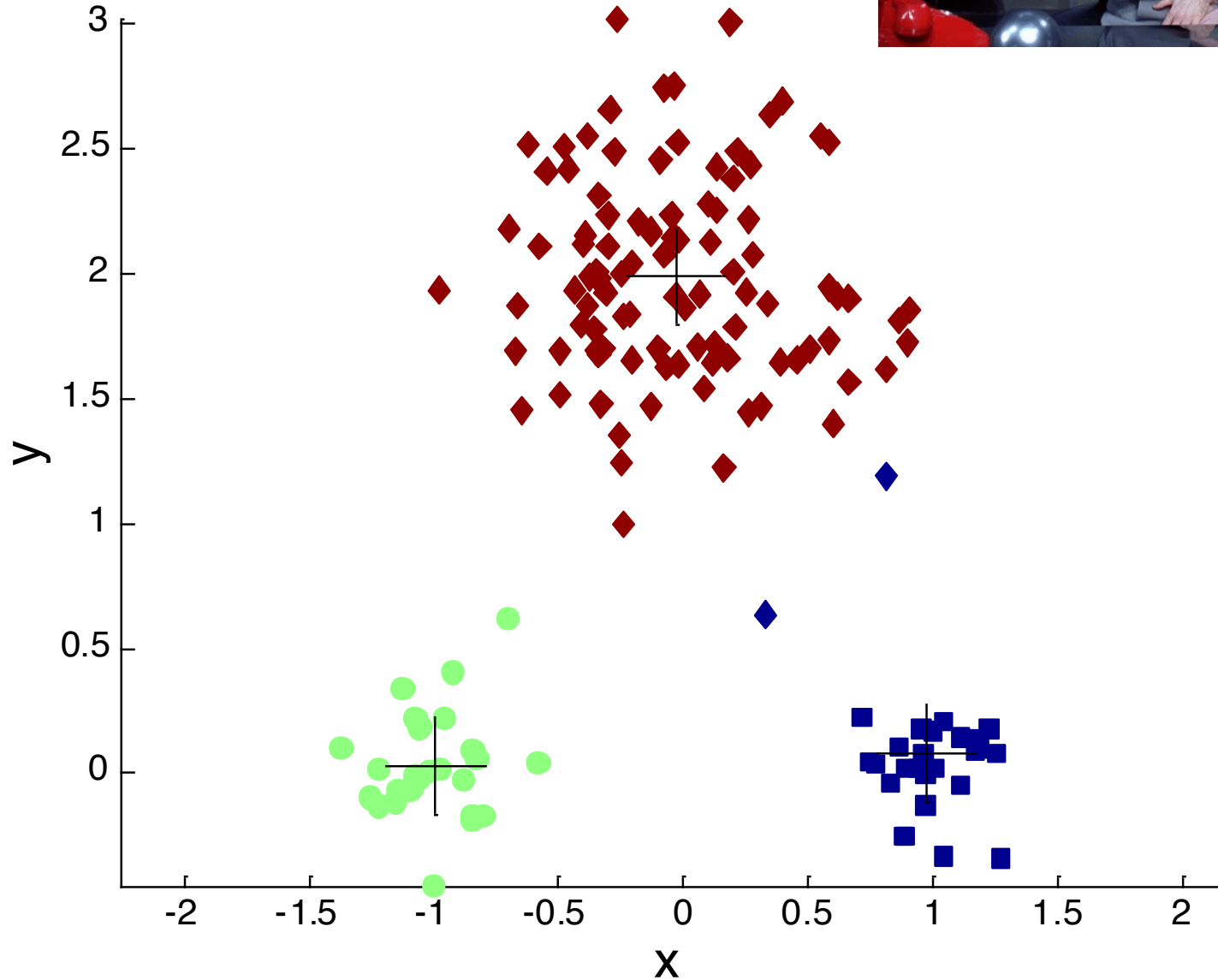


**Original Points**

**K-means Clusters**



*Oh, this might be easier than I expected.*
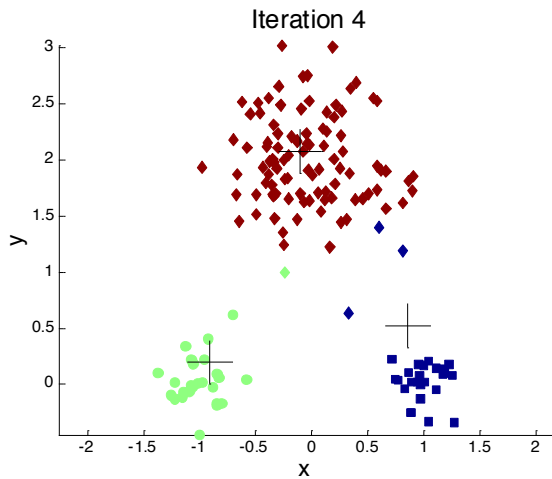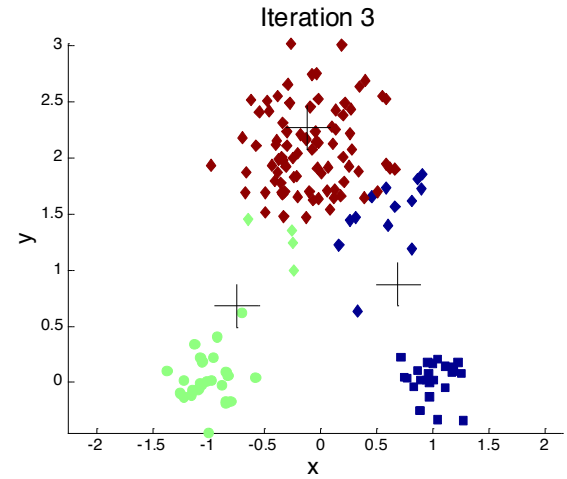
# Hold on... Another issue



Iteration 6

# Importance of Choosing Initial Centroids

# Importance of Choosing Initial Centroids ...

# Importance of Choosing Initial Centroids ...

# 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
- Post-processing
- Generate a larger number of clusters and then perform a hierarchical clustering
- K-means variants e.g. bisecting K-means

# 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)

# Agglomerative Clustering Algorithm

- Most 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**

p1   p2   p3   p4   . . .   p9   p10   p11   p12

# Intermediate Situation

- After some merging steps, we have some clusters



|    | C1 | C2 | C3 | C4 | C5 |
|----|----|----|----|----|----|
| C1 |    |    |    |    |    |
| C2 |    |    |    |    |    |
| C3 |    |    |    |    |    |
| C4 |    |    |    |    |    |
| C5 |    |    |    |    |    |

**Proximity Matrix**

p1  p2  p3  p4  ...  p9  p10  p11  p12

# 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**

C3

C4

C1

C2    C5

p1   p2   p3   p4   p9   p10   p11   p12

# After Merging

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

|          | C1 | C2 ∪ C5 | C3 | C4 |
|----------|----|---------|----|----|
| **C1**   |    | ?       |    |    |
| **C2 ∪ C5** | ? | ?    | ?  | ?  |
| **C3**   |    | ?       |    |    |
| **C4**   |    | ?       |    |    |

**Proximity Matrix**

# How to Define Inter-Cluster Distance

**Similarity?**

|     | p1 | p2 | p3 | p4 | p5 | . . . |
|-----|----|----|----|----|----|-------|
| **p1** |    |    |    |    |    |       |
| **p2** |    |    |    |    |    |       |
| **p3** |    |    |    |    |    |       |
| **p4** |    |    |    |    |    |       |
| **p5** |    |    |    |    |    |       |
| .   |    |    |    |    |    |       |
| .   |    |    |    |    |    |       |
| .   |    |    |    |    |    |       |

**Proximity Matrix**

- 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

| | p1 | p2 | p3 | p4 | p5 | . . . |
|-----|----|----|----|----|----|-------|
| p1 | | | | | | |
| p2 | | | | | | |
| p3 | | | | | | |
| p4 | | | | | | |
| p5 | | | | | | |
| . | | | | | | |

**Proximity Matrix**
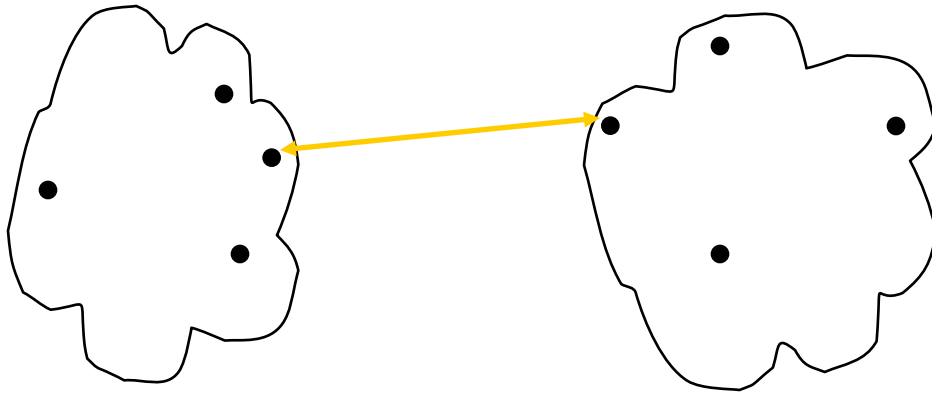
- <span style="color:red">MIN</span>
- MAX
- Group Average
- Distance Between Centroids
- Other methods driven by an objective function
  - Ward's Method uses squared error

# How to Define Inter-Cluster Similarity



|     | p1  | p2  | p3  | p4  | p5  | . . . |
|-----|-----|-----|-----|-----|-----|-------|
| p1  |     |     |     |     |     |       |
| p2  |     |     |     |     |     |       |
| p3  |     |     |     |     |     |       |
| p4  |     |     |     |     |     |       |
| p5  |     |     |     |     |     |       |
| .   |     |     |     |     |     |       |
| .   |     |     |     |     |     |       |
| .   |     |     |     |     |     |       |

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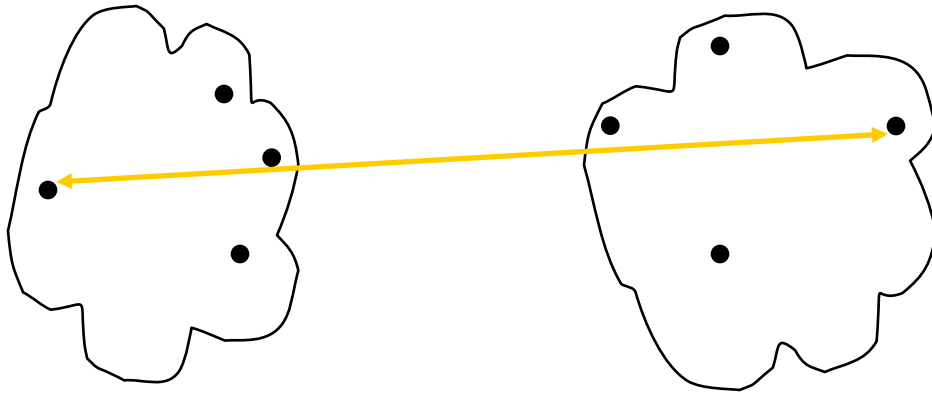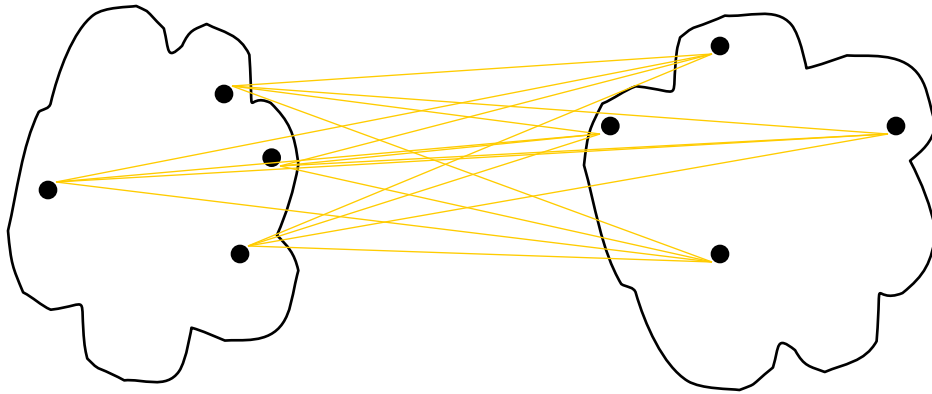
- MIN
- MAX
- Group Average
- Distance Between Centroids
- Other methods driven by an objective function
  - Ward's Method uses squared error

# How to Define Inter-Cluster Similarity



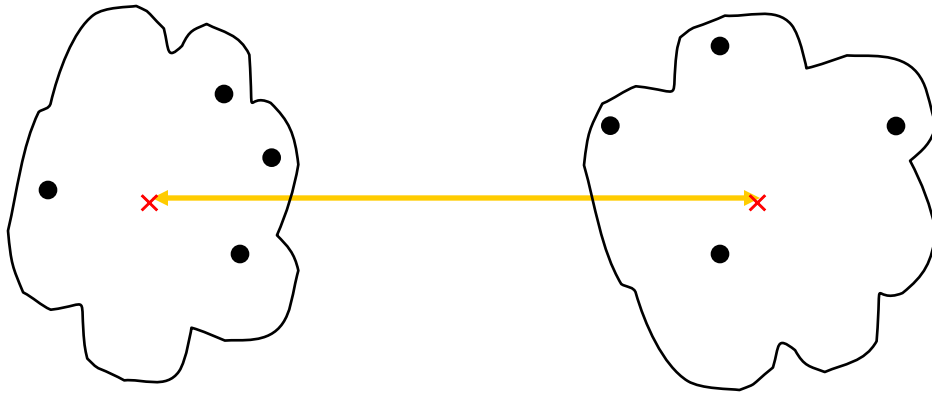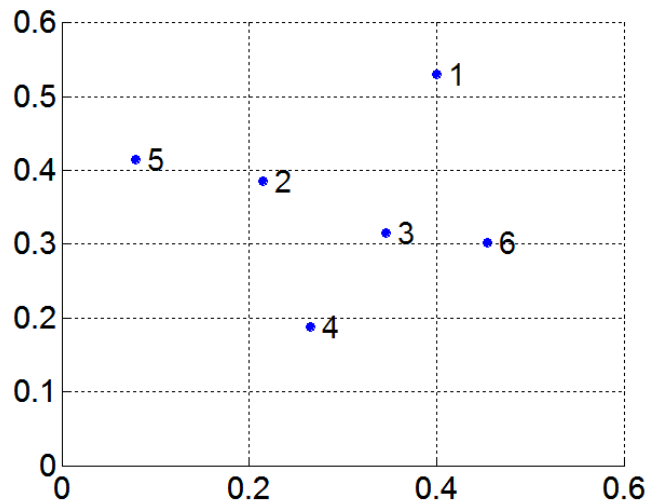| | p1 | p2 | p3 | p4 | p5 | . . . |
|------|----|----|----|----|----|-------|
| **p1** | | | | | | |
| **p2** | | | | | | |
| **p3** | | | | | | |
| **p4** | | | | | | |
| **p5** | | | | | | |
| **.** | | | | | | |
| **.** | | | | | | |
| **.** | | | | | | |

**Proximity Matrix**

- 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



|     | p1 | p2 | p3 | p4 | p5 | . . . |
|-----|----|----|----|----|----|-------|
| **p1** |    |    |    |    |    |       |
| **p2** |    |    |    |    |    |       |
| **p3** |    |    |    |    |    |       |
| **p4** |    |    |    |    |    |       |
| **p5** |    |    |    |    |    |       |
| **.** |    |    |    |    |    |       |
| **.** |    |    |    |    |    |       |
| **.** |    |    |    |    |    |       |

**Proximity Matrix**

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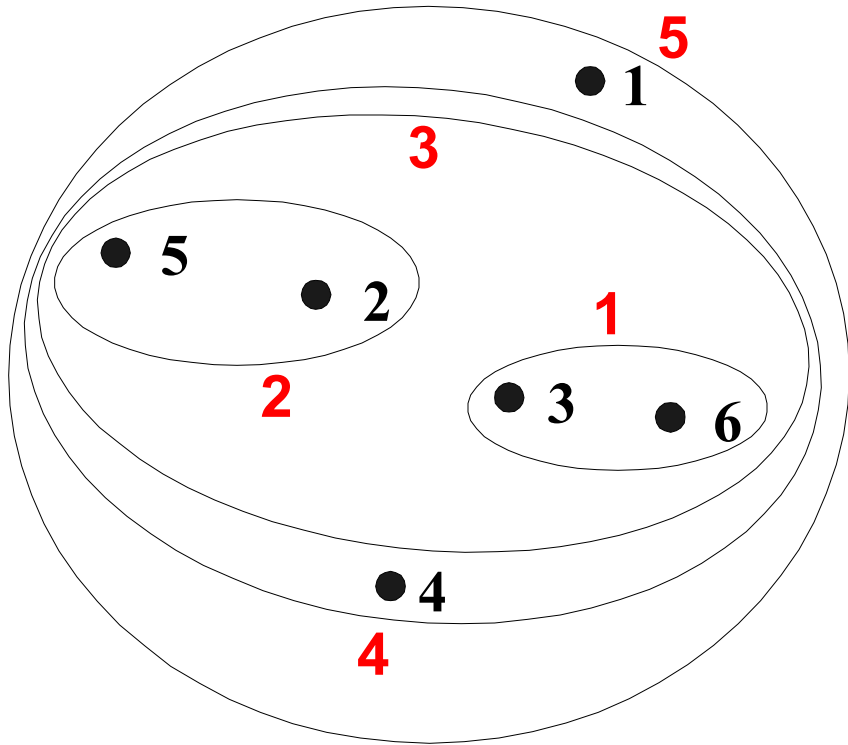  - Ward's Method uses squared error

# MIN or Single Link

- Proximity of two clusters is based on the two closest points in the different clusters
  - Determined by one pair of points, i.e., by one link in the proximity graph
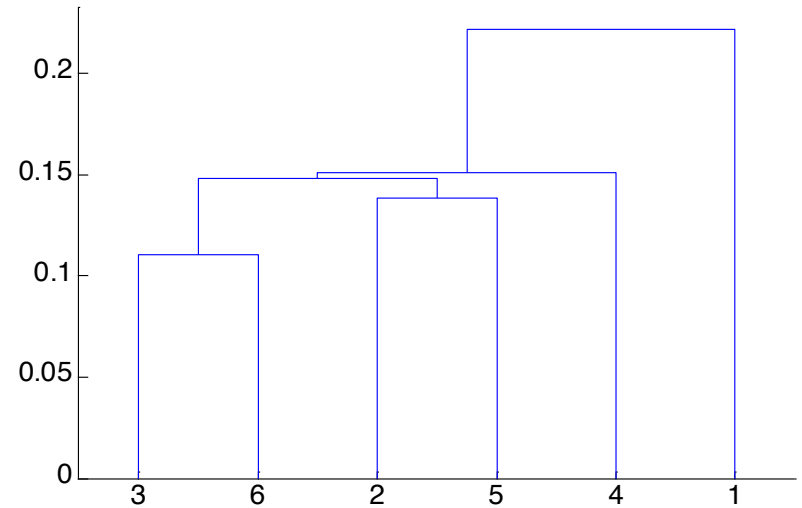
- Example:

**Distance Matrix:**

|     | p1   | p2   | p3   | p4   | p5   | p6   |
|-----|------|------|------|------|------|------|
| p1  | 0.00 | 0.24 | 0.22 | 0.37 | 0.34 | 0.23 |
| p2  | 0.24 | 0.00 | 0.15 | 0.20 | 0.14 | 0.25 |
| p3  | 0.22 | 0.15 | 0.00 | 0.15 | 0.28 | 0.11 |
| p4  | 0.37 | 0.20 | 0.15 | 0.00 | 0.29 | 0.22 |
| p5  | 0.34 | 0.14 | 0.28 | 0.29 | 0.00 | 0.39 |
| p6  | 0.23 | 0.25 | 0.11 | 0.22 | 0.39 | 0.00 |

# Hierarchical Clustering: MIN
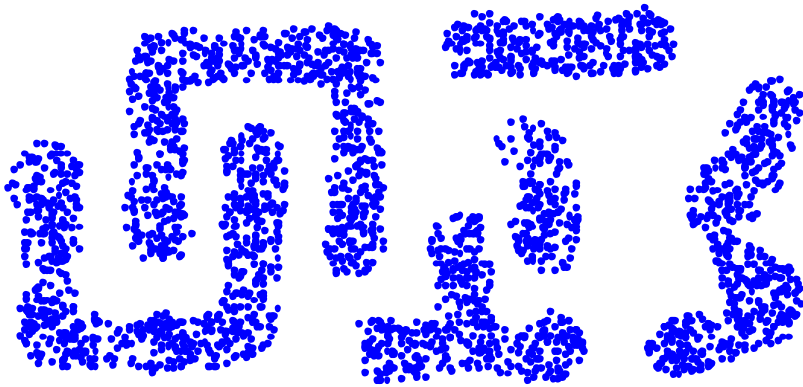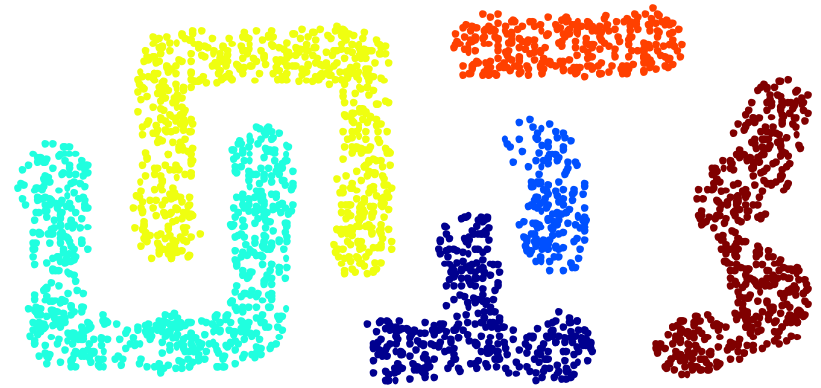


**Nested Clusters**

**Dendrogram**

# Strength of MIN
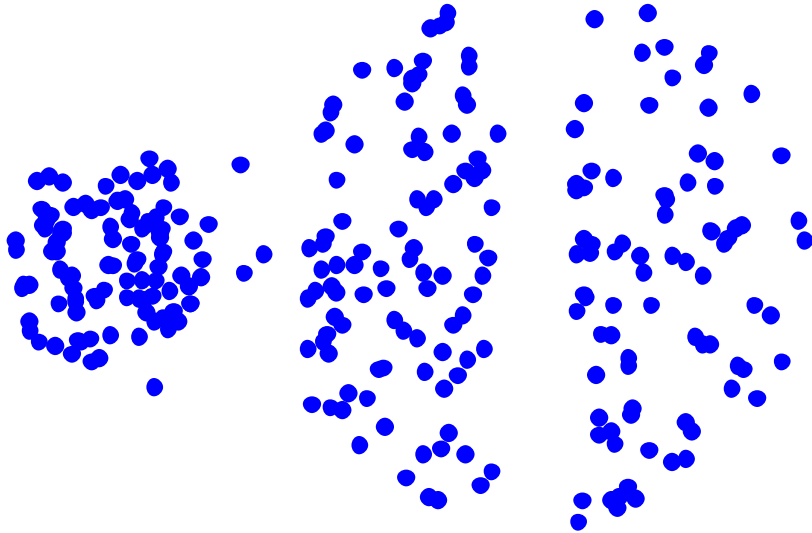


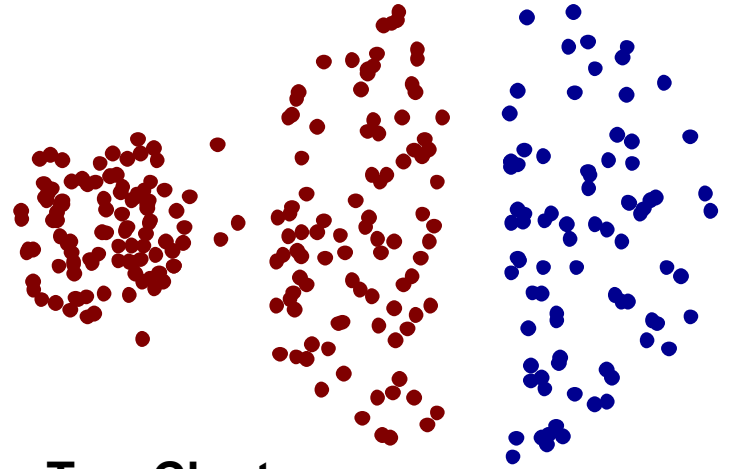**Original Points**                    **Six Clusters**

- Can handle non-elliptical shapes
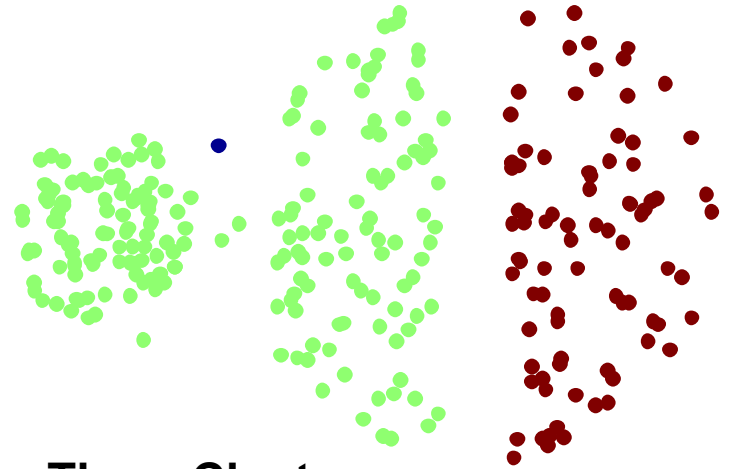
# Limitations of MIN



**Original Points**

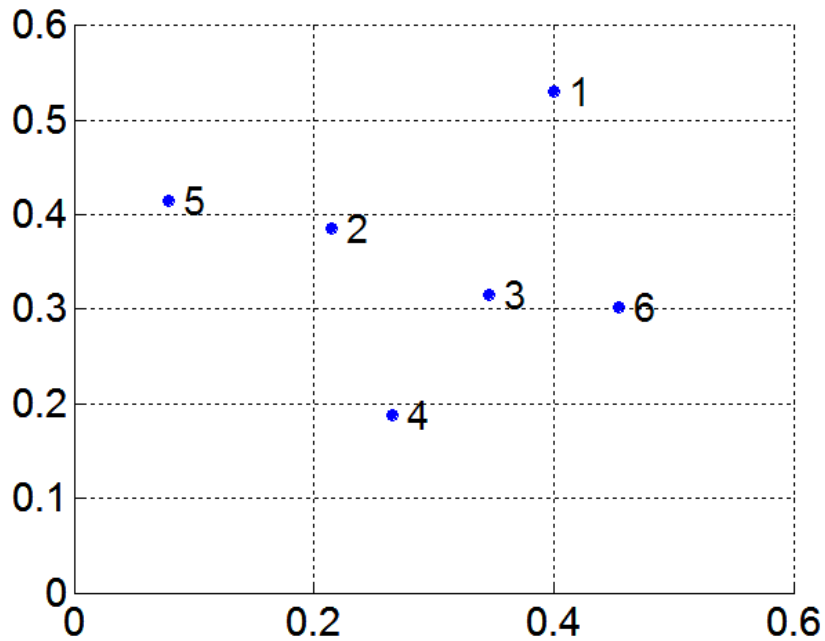- **Sensitive to noise and outliers**

**Two Clusters**

**Three Clusters**
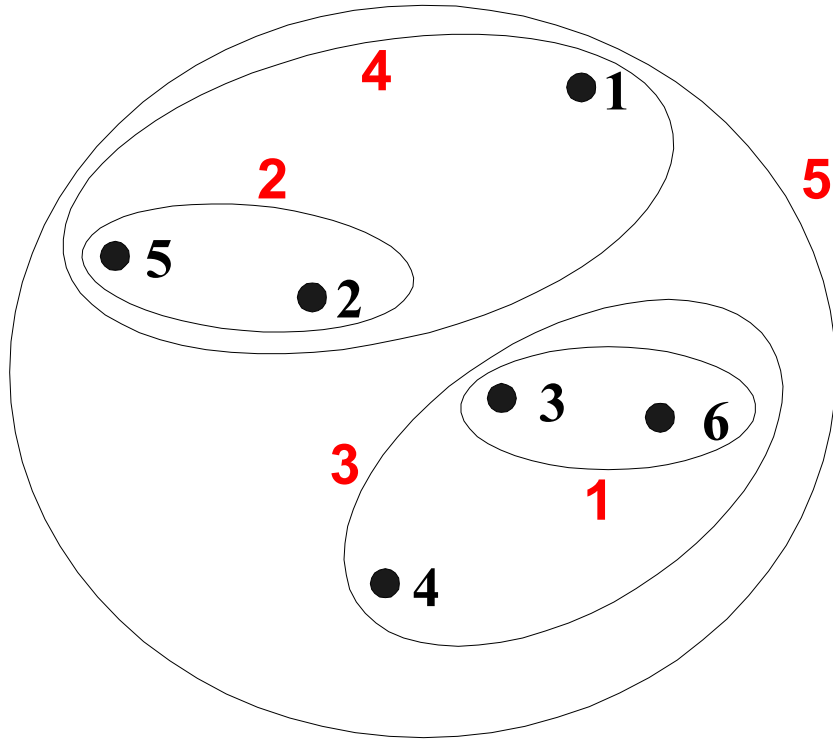
# MAX or Complete Linkage

- Proximity of two clusters is based on the two most distant points in the different clusters
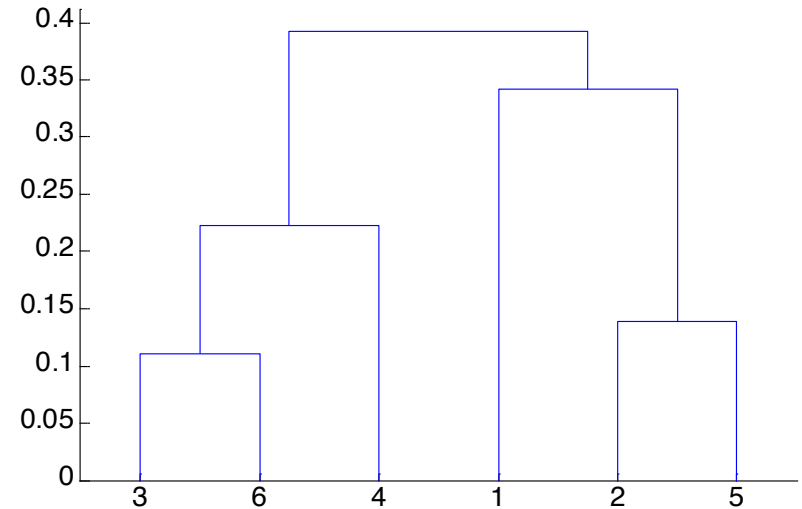  - Determined by all pairs of points in the two clusters



**Distance Matrix:**

|  | p1 | p2 | p3 | p4 | p5 | p6 |
|----|------|------|------|------|------|------|
| p1 | 0.00 | 0.24 | 0.22 | 0.37 | 0.34 | 0.23 |
| p2 | 0.24 | 0.00 | 0.15 | 0.20 | 0.14 | 0.25 |
| p3 | 0.22 | 0.15 | 0.00 | 0.15 | 0.28 | 0.11 |
| p4 | 0.37 | 0.20 | 0.15 | 0.00 | 0.29 | 0.22 |
| p5 | 0.34 | 0.14 | 0.28 | 0.29 | 0.00 | 0.39 |
| p6 | 0.23 | 0.25 | 0.11 | 0.22 | 0.39 | 0.00 |

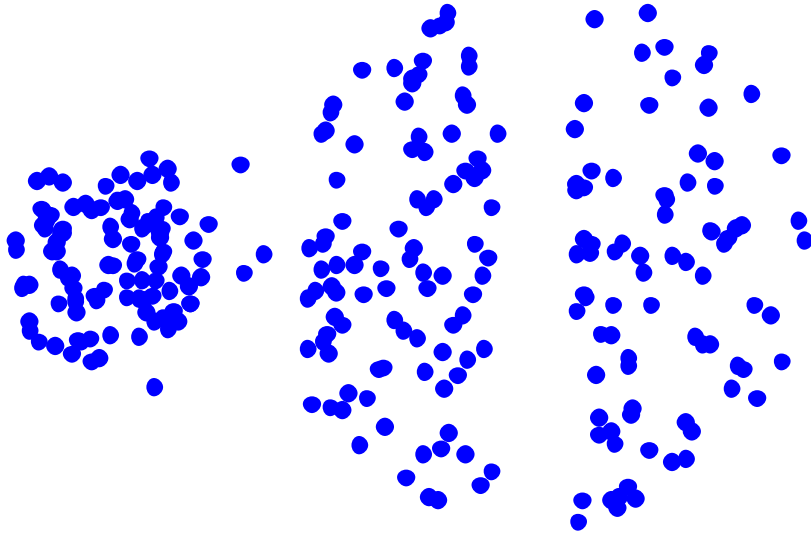# Hierarchical Clustering: MAX

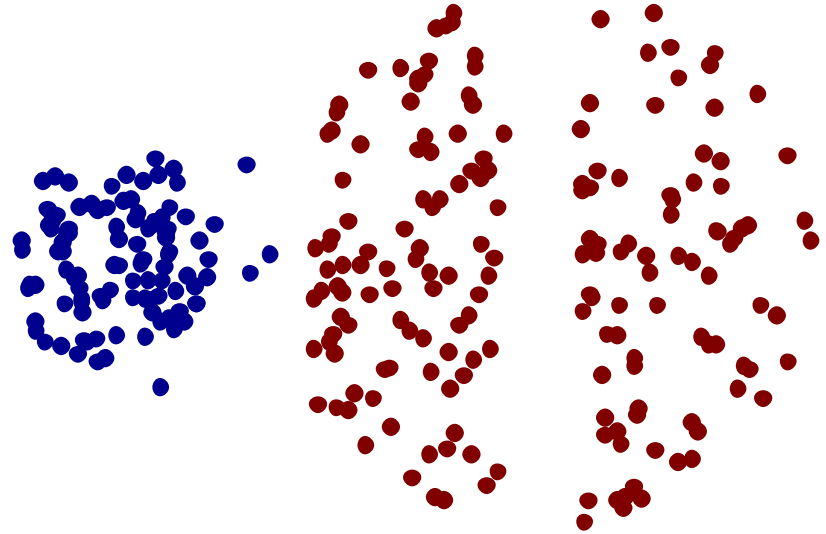

**Nested Clusters**

**Dendrogram**

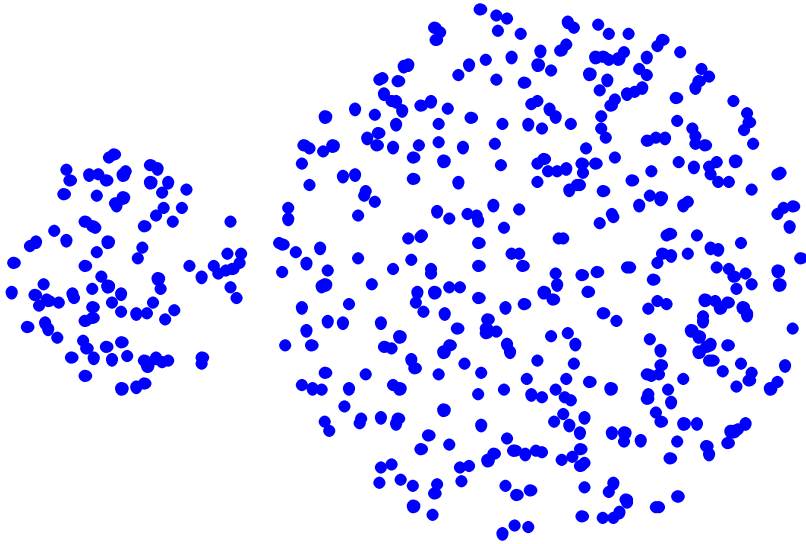# Strength of MAX



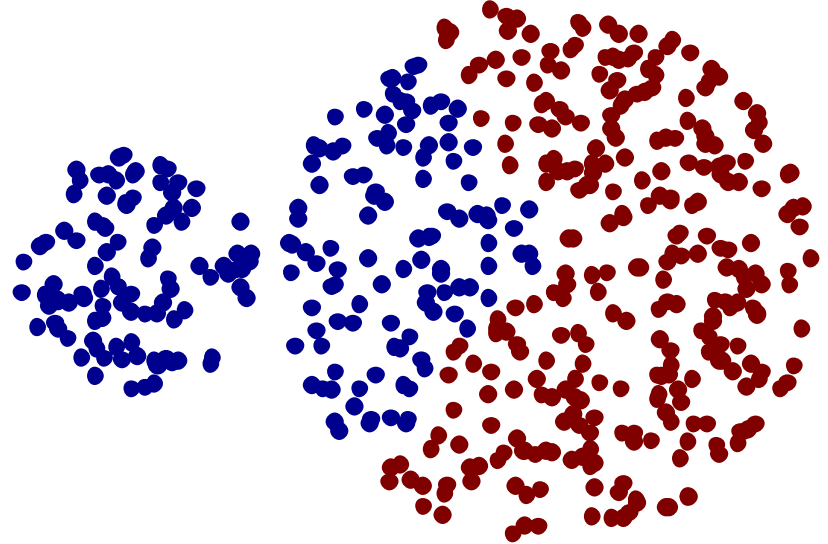**Original Points**

**Two Clusters**

- **Less susceptible to noise and outliers**

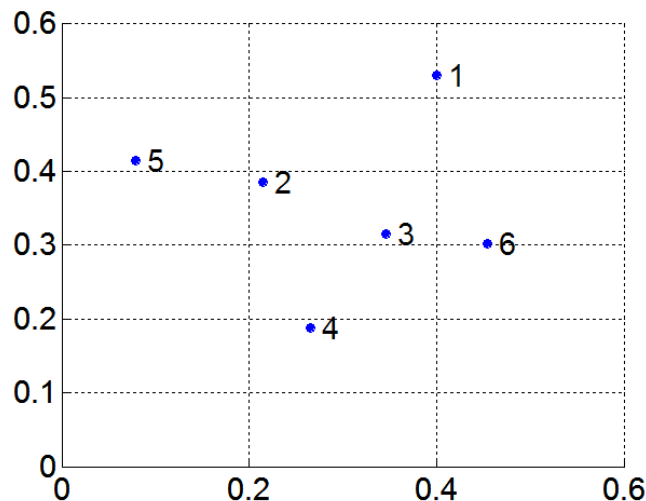# Limitations of MAX



**Original Points**

**Two Clusters**

- **Tends to break large clusters**
- **Biased towards globular clusters**

# Group Average

- Proximity of two clusters is the average of pairwise proximity between points in the two clusters.

$$\text{proximity}(\text{Cluster}_i, \text{Cluster}_j) = \frac{\sum_{\substack{p_i \in \text{Cluster}_i \\ p_j \in \text{Cluster}_j}} \text{proximity}(p_i, p_j)}{|\text{Cluster}_i| \times |\text{Cluster}_j|}$$
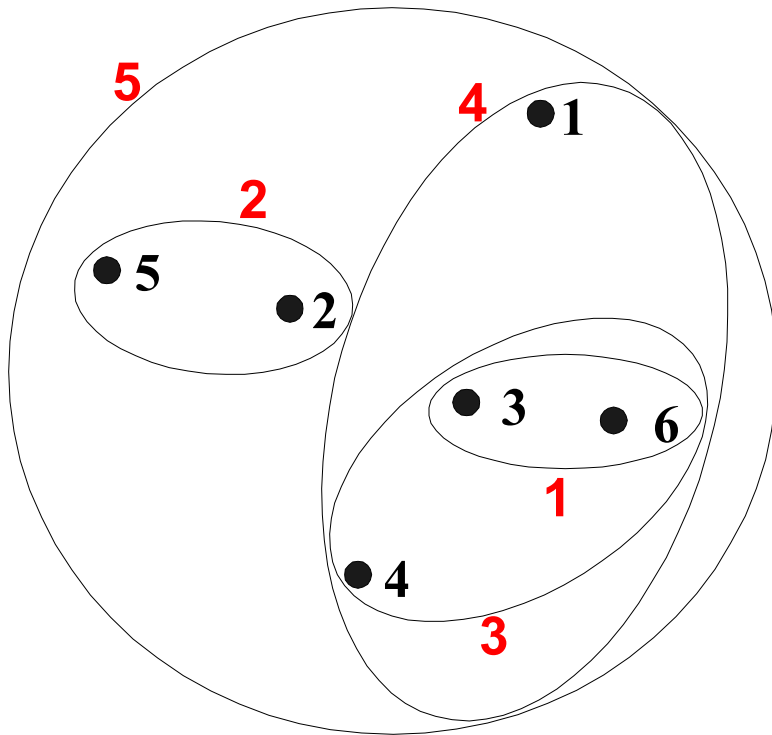
- Need to use average connectivity for scalability since total proximity favors large clusters
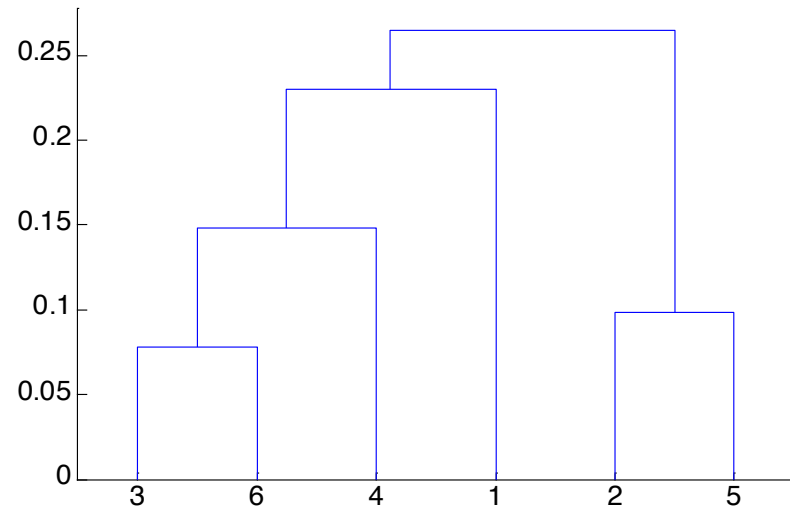
**Distance Matrix:**

|    | p1   | p2   | p3   | p4   | p5   | p6   |
|----|------|------|------|------|------|------|
| p1 | 0.00 | 0.24 | 0.22 | 0.37 | 0.34 | 0.23 |
| p2 | 0.24 | 0.00 | 0.15 | 0.20 | 0.14 | 0.25 |
| p3 | 0.22 | 0.15 | 0.00 | 0.15 | 0.28 | 0.11 |
| p4 | 0.37 | 0.20 | 0.15 | 0.00 | 0.29 | 0.22 |
| p5 | 0.34 | 0.14 | 0.28 | 0.29 | 0.00 | 0.39 |
| p6 | 0.23 | 0.25 | 0.11 | 0.22 | 0.39 | 0.00 |

# 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: Problems and Limitations

- Once a decision is made to combine two clusters, it cannot be undone

- It needs much more space AND time

- No global objective function is directly minimized

- Different schemes have problems with one or more of the following:

  - Sensitivity to noise and outliers

  - Difficulty handling clusters of different sizes and non-globular shapes
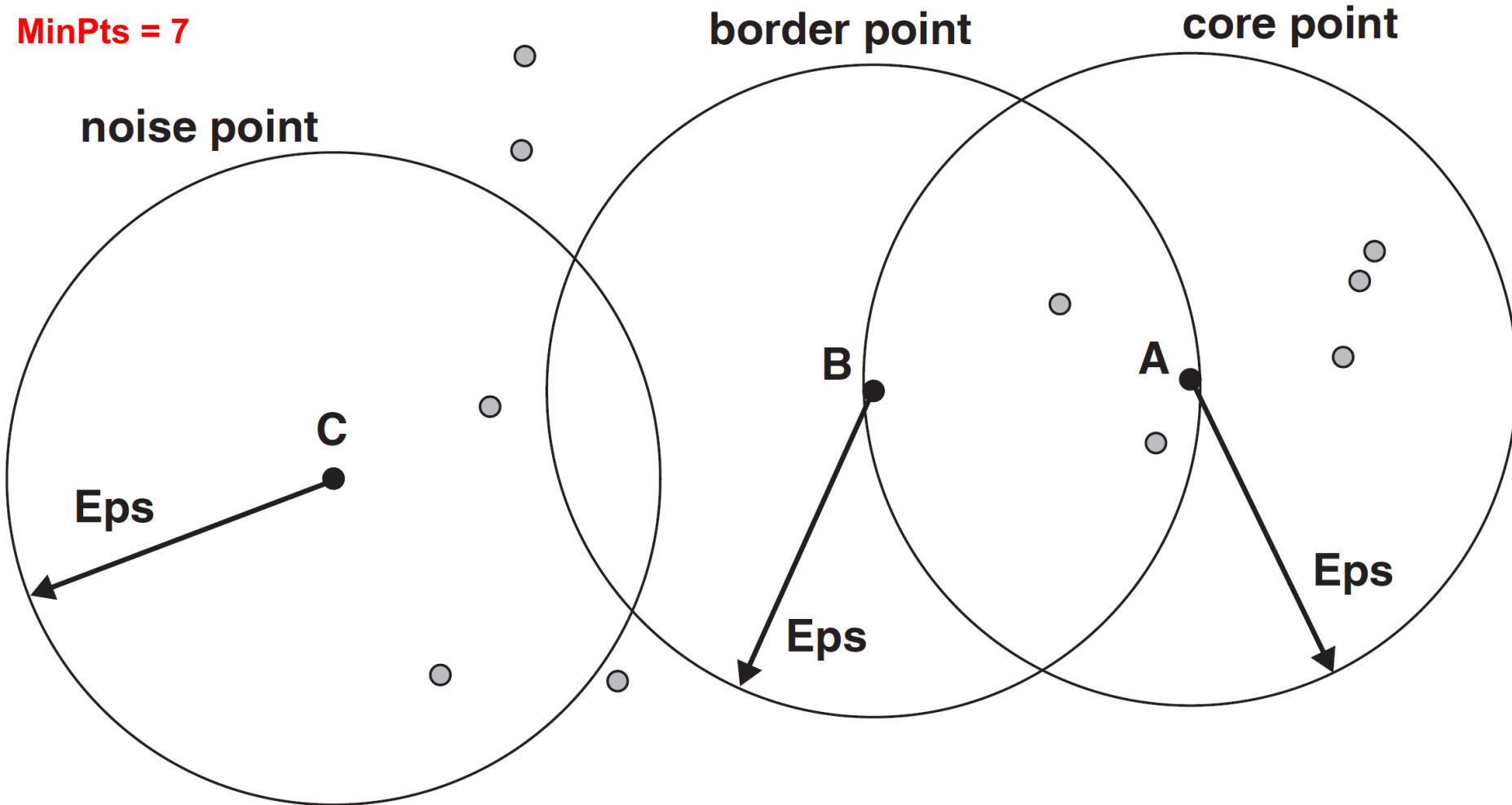
  - Breaking large clusters

# DBSCAN

- ## DBSCAN is a **density-based** algorithm.

  - Density = number of points within a specified radius (**Eps**)

  - A point is a core point if it has at least a specified number of points (**MinPts**) within Eps
    - These are points that are at the interior of a cluster
    - Counts the point itself

  - A border point is not a core point, but is in the neighborhood of a core point

  - A noise point is any point that is not a core point or a border point

# DBSCAN: Core, Border, Noise Points



MinPts = 7

border point

core point

noise point

B

A

C

Eps

Eps

Eps

# 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 $current\_cluster\_label$
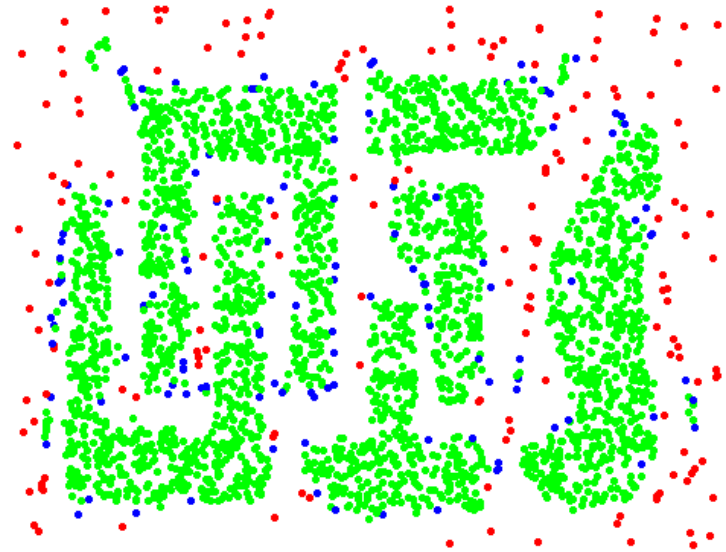
        **end if**

    **end for**

**end for**

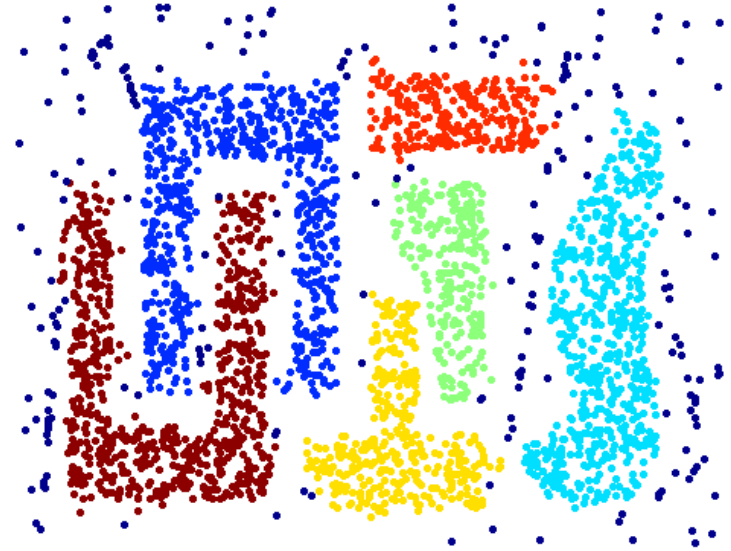# DBSCAN: Core, Border and Noise Points



**Original Points**

Point types: **<span style="color:green">core</span>**, **<span style="color:blue">border</span>** and **<span style="color:red">noise</span>**
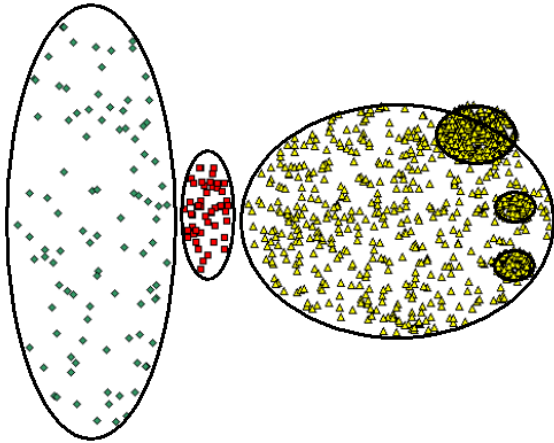
**Eps = 10, MinPts = 4**

# When DBSCAN Works Well
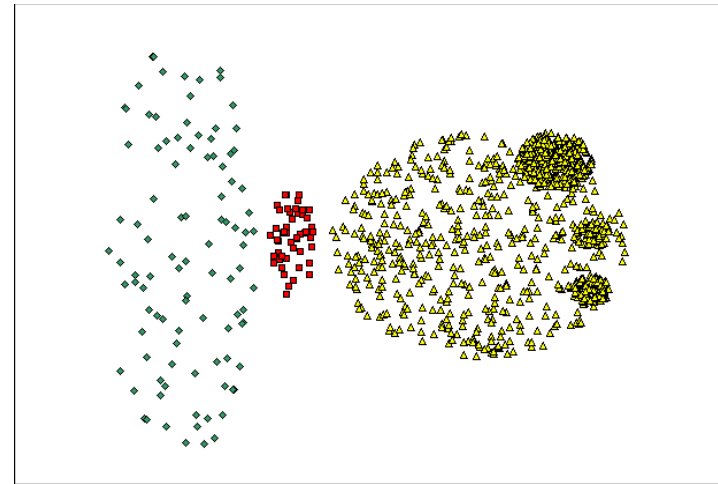


**Original Points**

**Clusters**

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

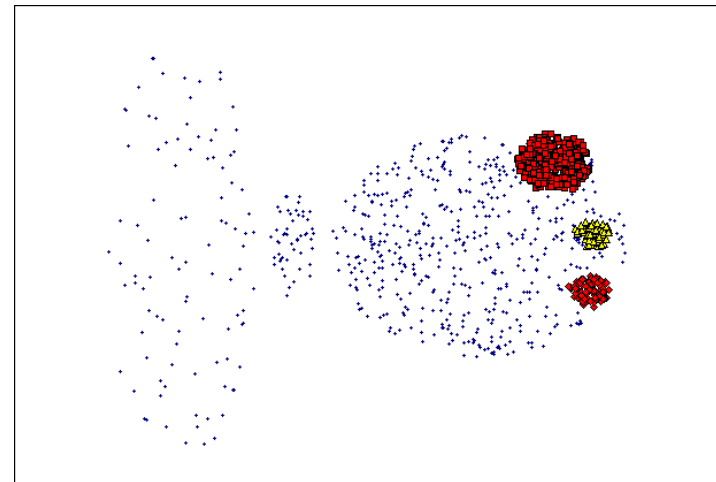# When DBSCAN Does NOT Work Well



**Original Points**
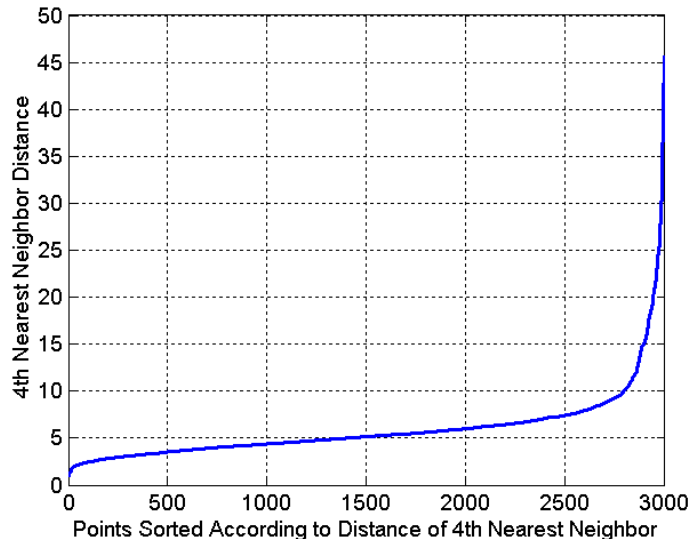
- **Varying densities**
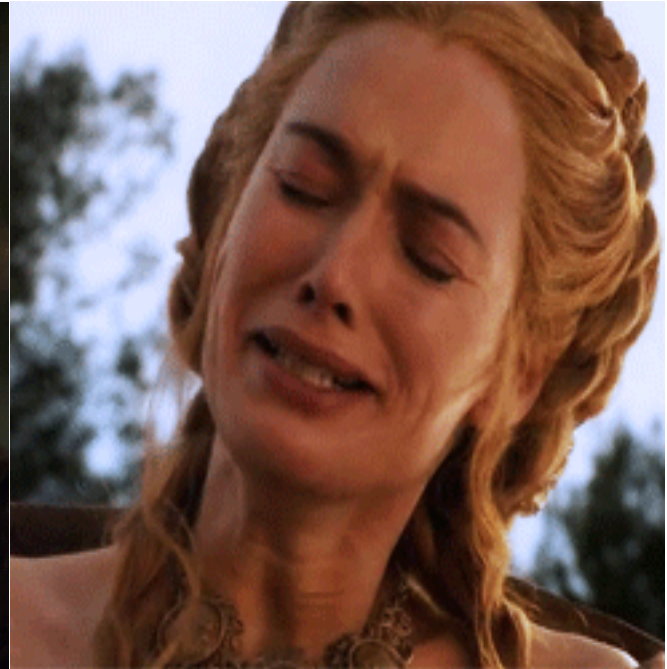
- **High-dimensional data**

(MinPts=4, Eps=9.75).

(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 further distance

- So, plot sorted distance of every point to its $k^{th}$ nearest neighbor
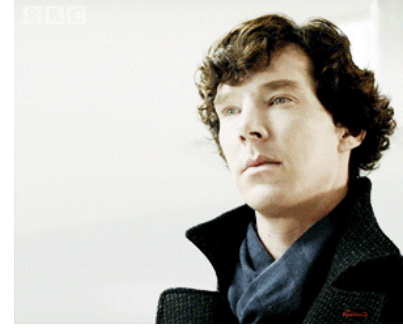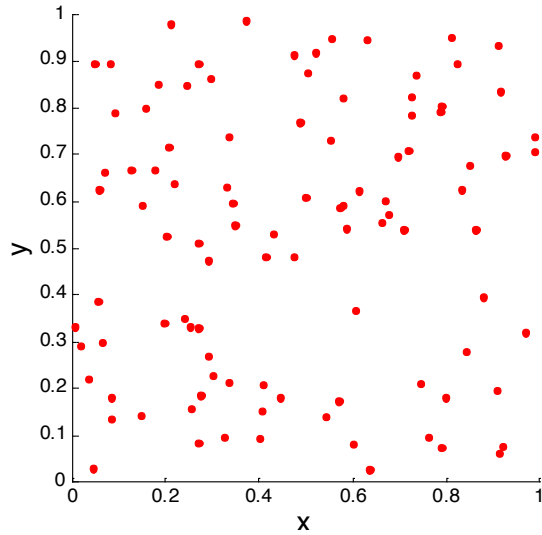
# Cluster Validity

# 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"!

- Evaluation is really important here:
  - To avoid finding patterns in noise
  - To compare clustering algorithms
  - To compare two sets of clusters
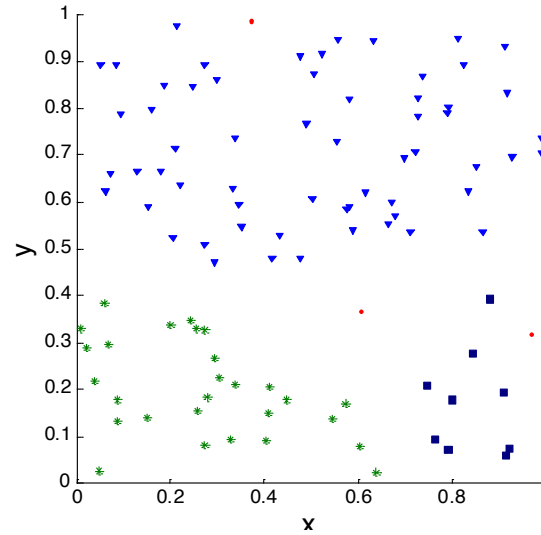  - To compare two clusters
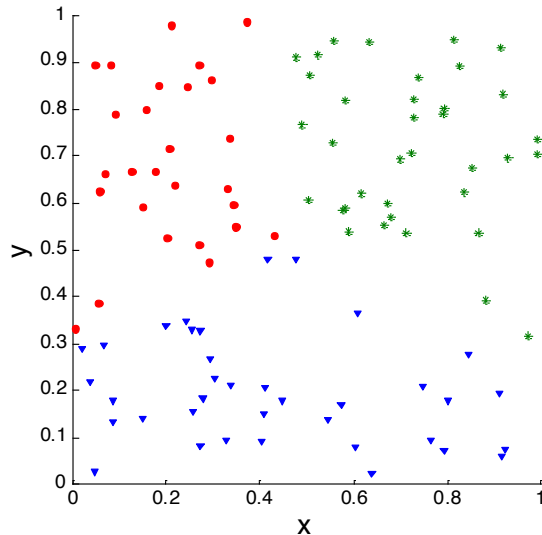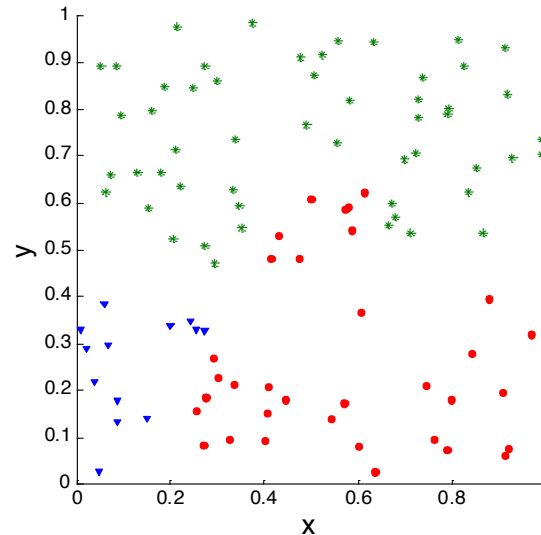
# Clusters found in Random Data



**Random Points**

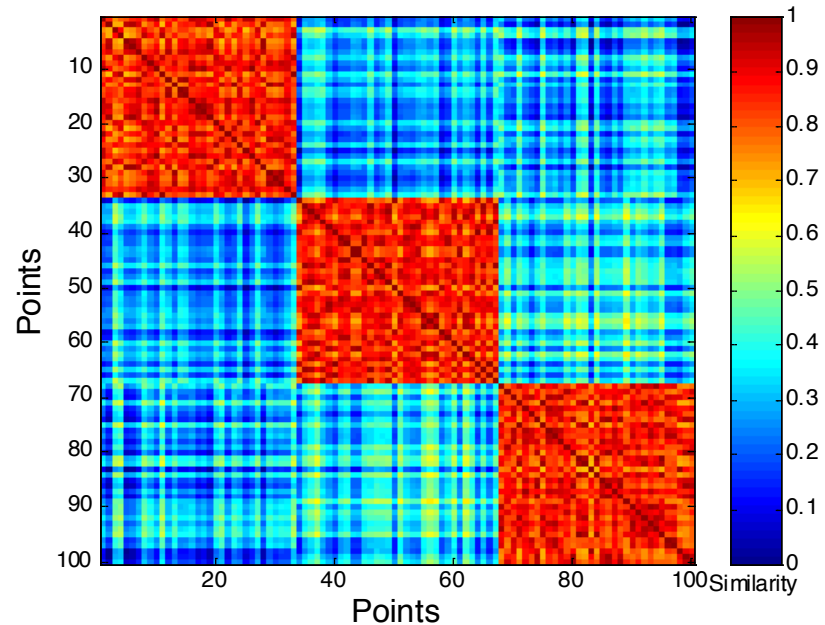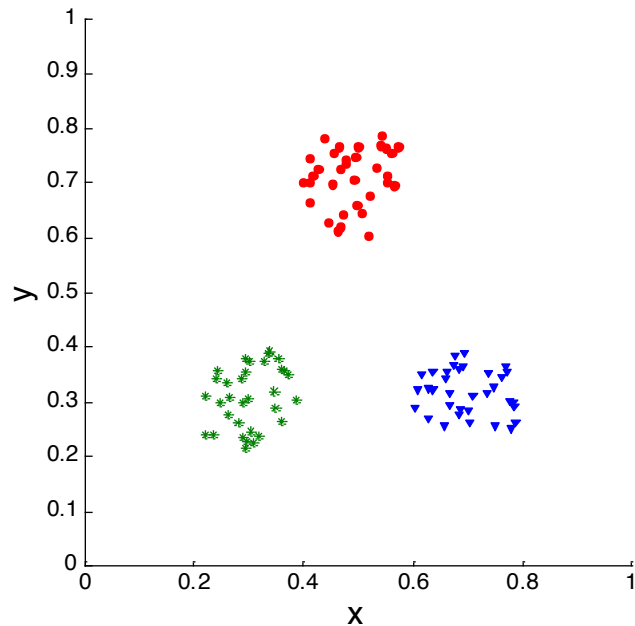**DBSCAN**

**K-means**

**Complete Link**

# 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 or Purity

  - 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.
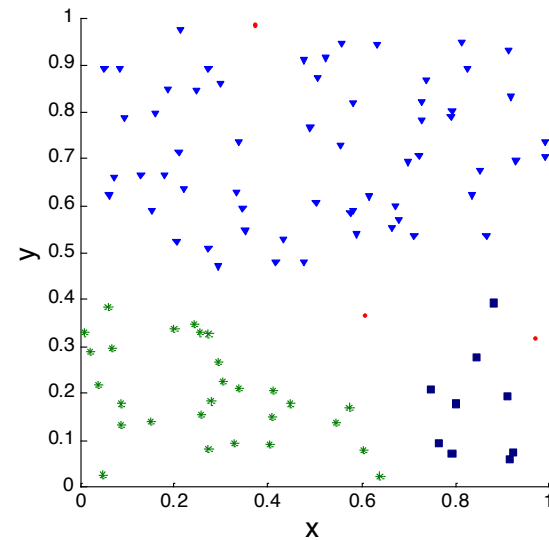    - ◆ Often an external or internal index is used for this function, e.g., SSE or entropy
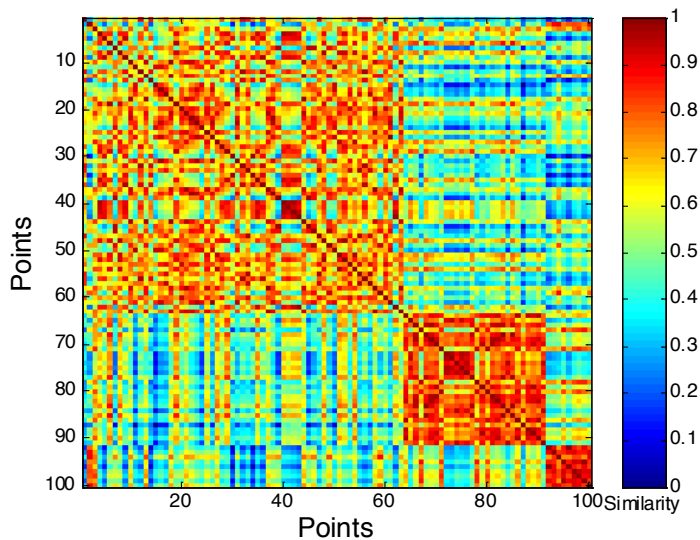
# Cluster Validity: Similarity Matrix

- **Order** the similarity matrix with respect to cluster labels and inspect visually.

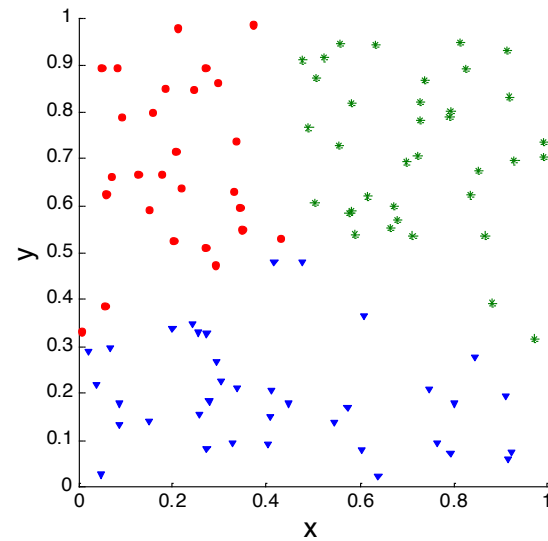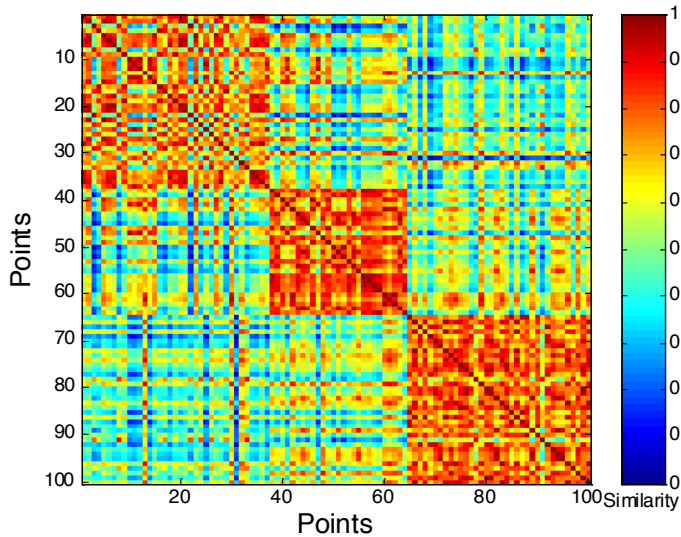# Cluster Validity: Similarity Matrix

- Clusters in random data are not so crisp
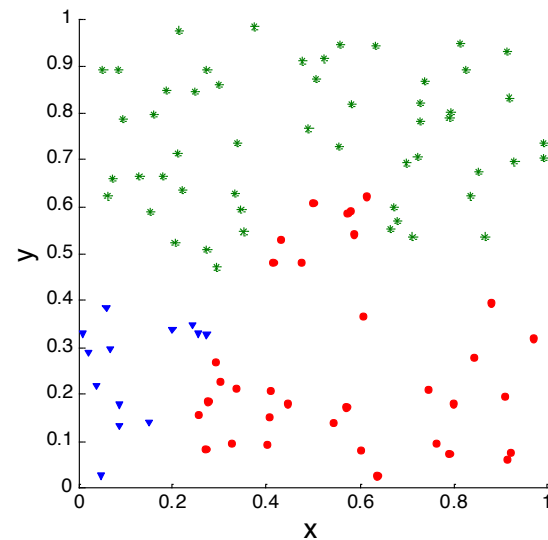


**DBSCAN**

# Cluster Validity: Similarity Matrix

- Clusters in random data are not so crisp



**K-means**

# Cluster Validity: Similarity Matrix

- Clusters in random data are not so crisp



**Complete Link**

# Cluster Validity: Similarity Matrix



**DBSCAN**

# Internal Measures: SSE

- Clusters in more complicated figures aren't well separated

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

  - SSE

- 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**

# 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 all links within a cluster.

   – Cluster separation is the sum of the weights between nodes in the cluster and nodes outside the cluster.

cohesion                                        separation

# Internal Measures: Silhouette Coefficient

- Silhouette coefficient combines 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 = (b − a) / max(a,b)

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



Distances used to calculate **b**

$i$

Distances used to calculate **a**

- Can calculate the average silhouette coefficient 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., $e = \sum_{i=1}^{K} \frac{m_i}{m} e_j$, 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 $purity_j = \max p_{ij}$ and the overall purity of a clustering by $purity = \sum_{i=1}^{K} \frac{m_i}{m} purity_j$.

# 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

# To summarize

- Clustering is the most basic unsupervised technique
- Different algorithms might raise different results for what is the "optimal" clustering
- It is important to properly evaluate the results and justify any conclusion/decision using numbers



Consider yourself warned.

May the Force be with you.