

# Introduction to Machine Learning Applications

Spring 2023

Clustering

**Minor Gordon**

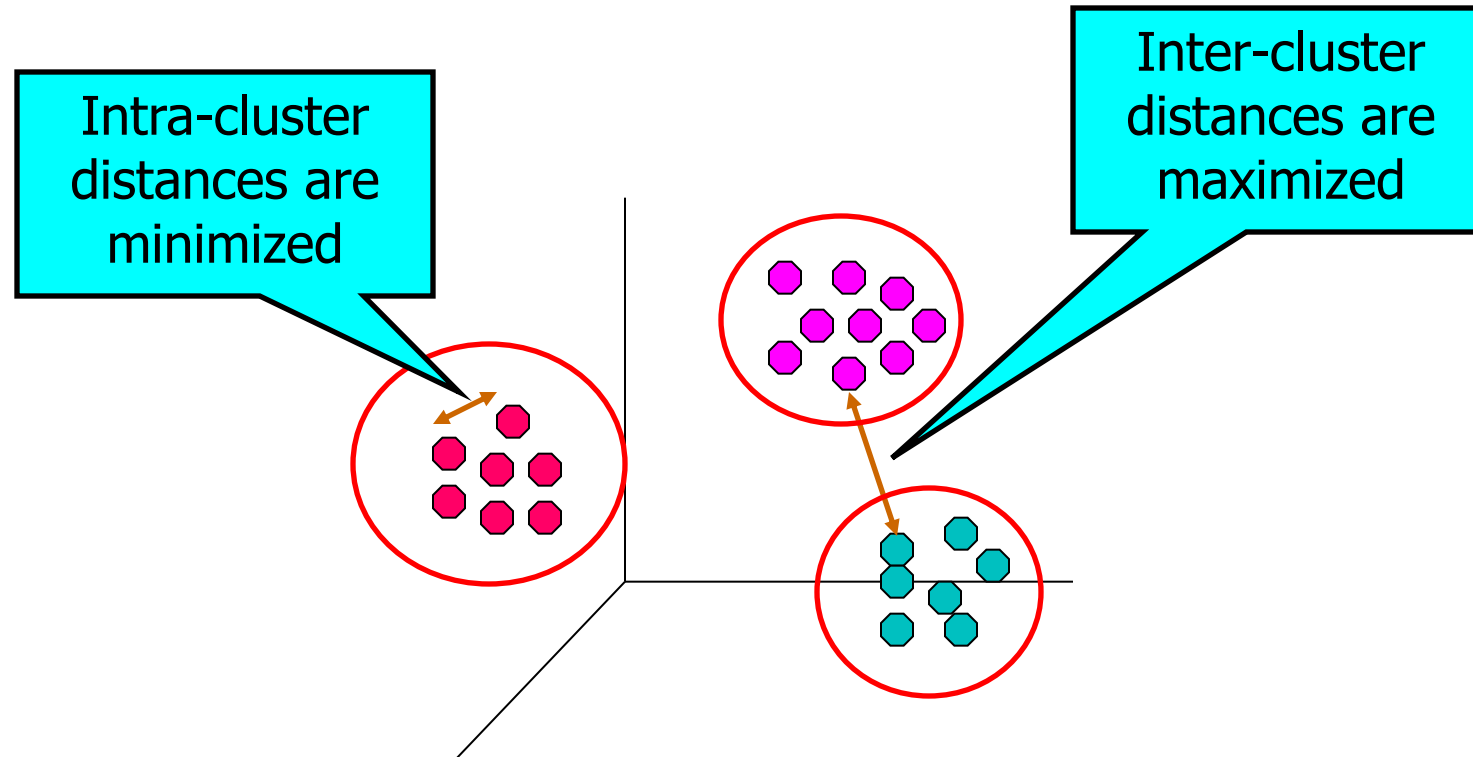
[gordom6@rpi.edu](mailto:gordom6@rpi.edu)



**Rensselaer**

# What is Cluster Analysis?

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



# Applications of Cluster Analysis

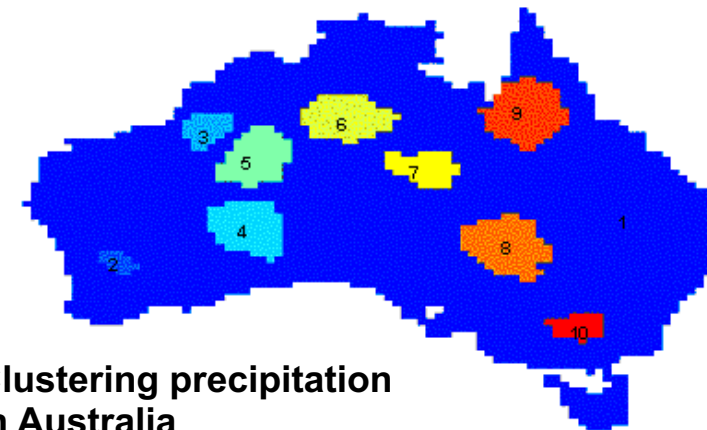
## ● Understanding

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

## ● Summarization

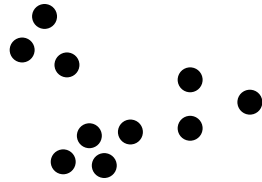
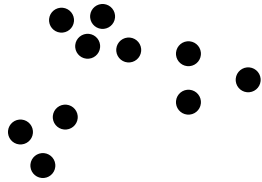
- Reduce the size of large data sets

	<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

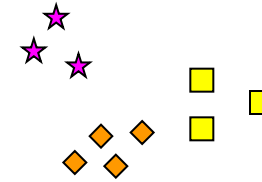
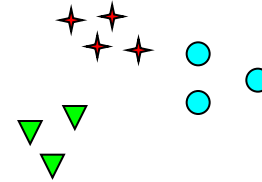


Clustering precipitation  
in Australia

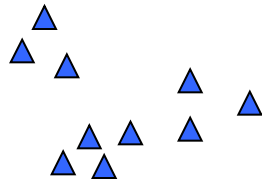
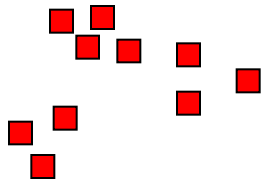
# Notion of a Cluster can be Ambiguous



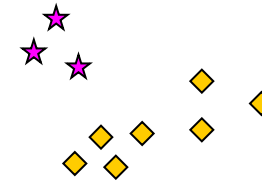
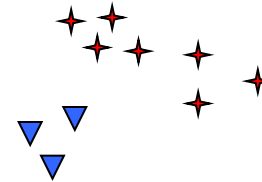
How many clusters?



Six Clusters



Two Clusters



Four Clusters

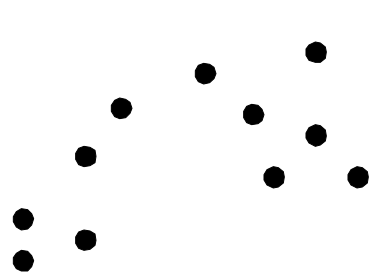
# Types of Clusterings

---

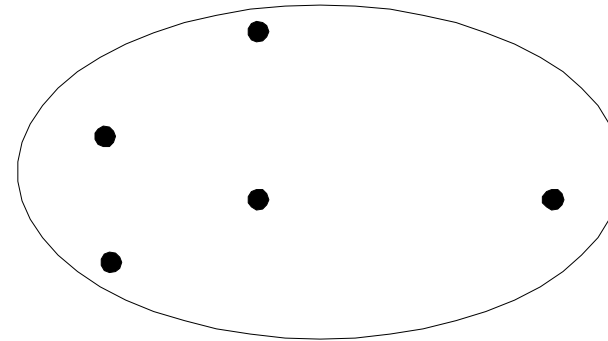
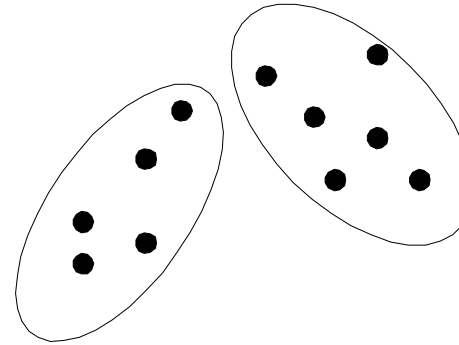
- A **clustering** is a set of clusters
- Important distinction between **hierarchical** and **partitional** sets of clusters
  - Partitional Clustering
    - ◆ A division of data objects into non-overlapping subsets (clusters)
  - Hierarchical clustering
    - ◆ A set of nested clusters organized as a hierarchical tree

# Partitional Clustering

---

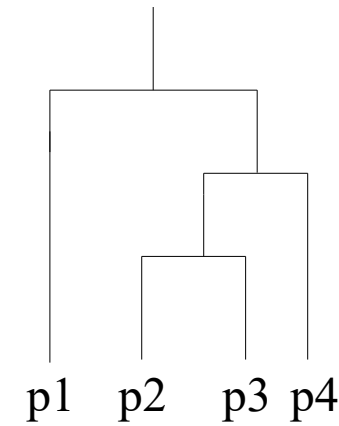
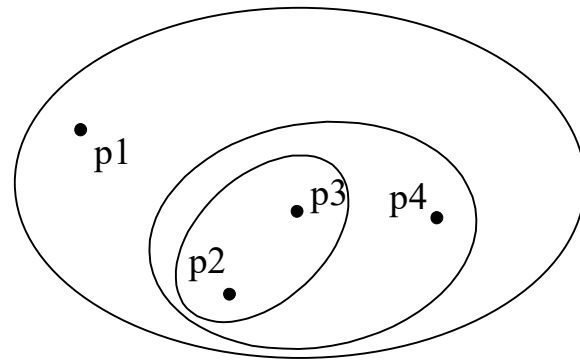


**Original Points**



**A Partitional Clustering**

# Hierarchical Clustering



# Other Distinctions Between Sets of Clusters

---

- Exclusive versus non-exclusive
  - In non-exclusive clusterings, points may belong to multiple clusters.
    - ◆ Can belong to multiple classes or could be 'border' points
  - Fuzzy clustering (one type of non-exclusive)
    - ◆ 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



# Types of Clusters

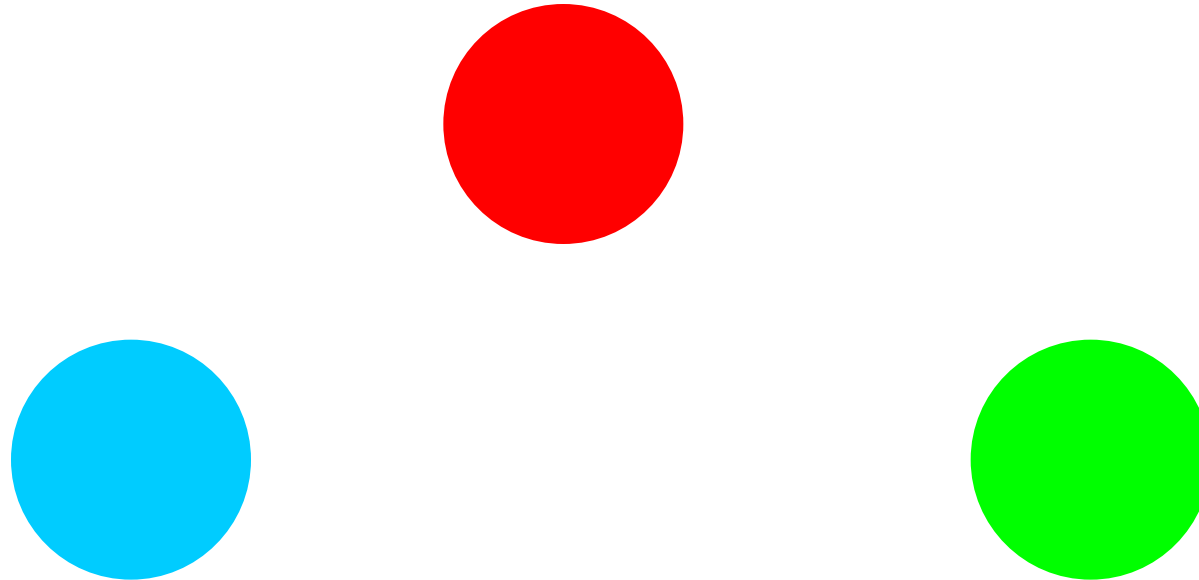
---

- Well-separated clusters
- Prototype-based clusters
- Contiguity-based clusters
- Density-based clusters
- 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.



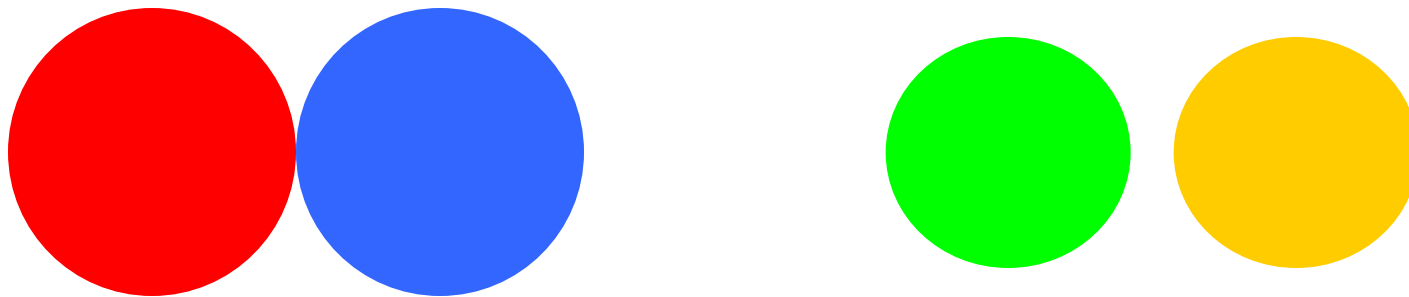
**3 well-separated clusters**

# Types of Clusters: Prototype-Based

---

- Prototype-based

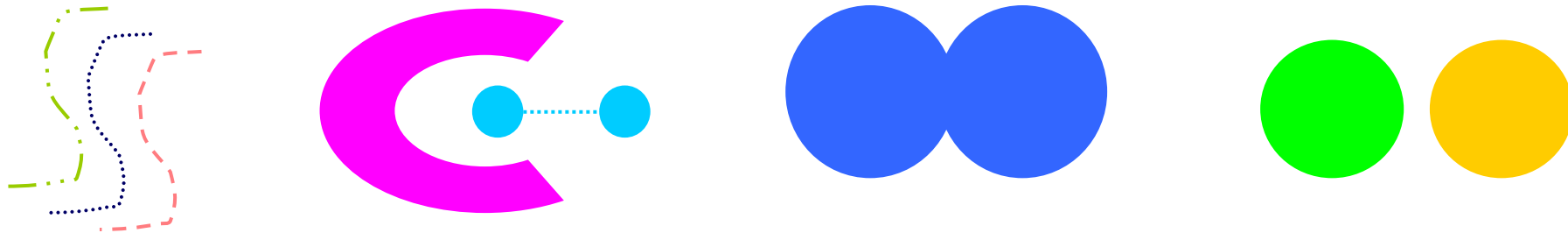
- A cluster is a set of objects such that an object in a cluster is closer (more similar) to the prototype or “center” of a cluster, than to the center of any other cluster
- The center of a cluster is often a **centroid**, the average of all the points in the cluster, or a **medoid**, the most “representative”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



**4 center-based clusters**

# Types of Clusters: Contiguity-Based

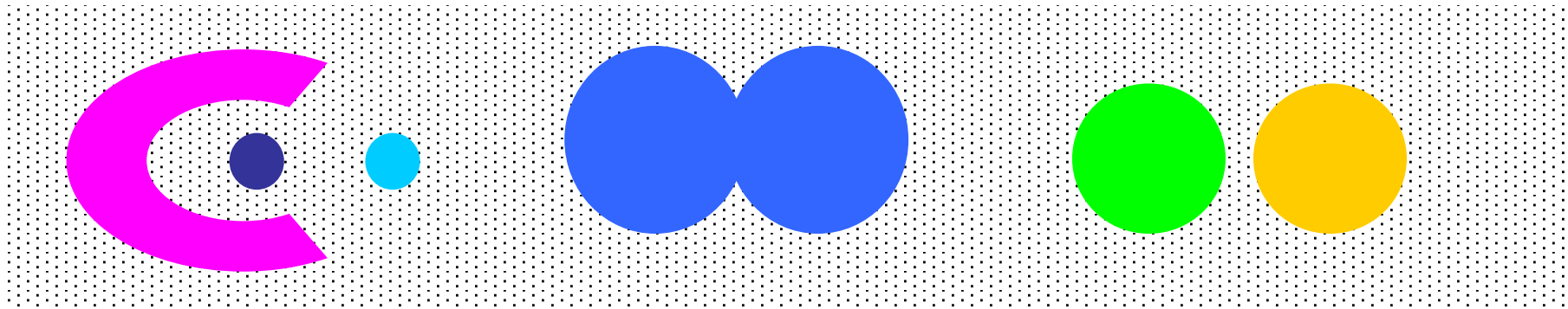
- Contiguous Cluster (Nearest neighbor or Transitive)
  - A cluster is a set of points such that a point in a cluster is closer (or more similar) to one or more other points in the cluster than to any point not in the cluster.



**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**

# Characteristics of the Input Data Are Important

---

- Type of proximity or density measure
  - Central to clustering
  - Depends on data and application
- Data characteristics that affect proximity and/or density are
  - Dimensionality
    - ◆ Sparseness
  - Attribute type
  - Special relationships in the data
    - ◆ For example, autocorrelation
  - Distribution of the data
- Noise and Outliers
  - Often interfere with the operation of the clustering algorithm
- Clusters of differing sizes, densities, and shapes

# Clustering Algorithms

---

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

# K-means Clustering

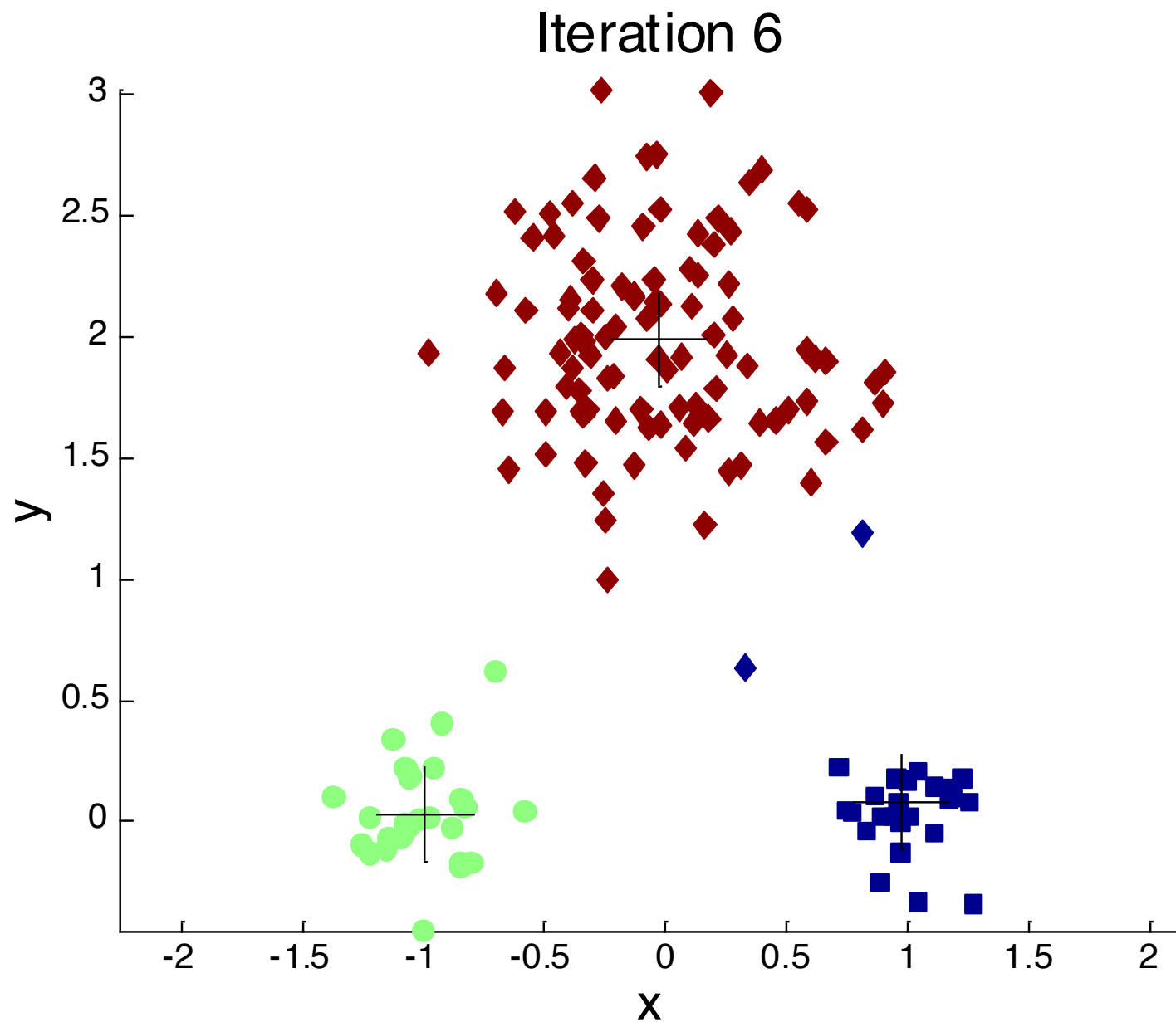
---

- Partitional clustering approach
- 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
- The basic algorithm is very simple

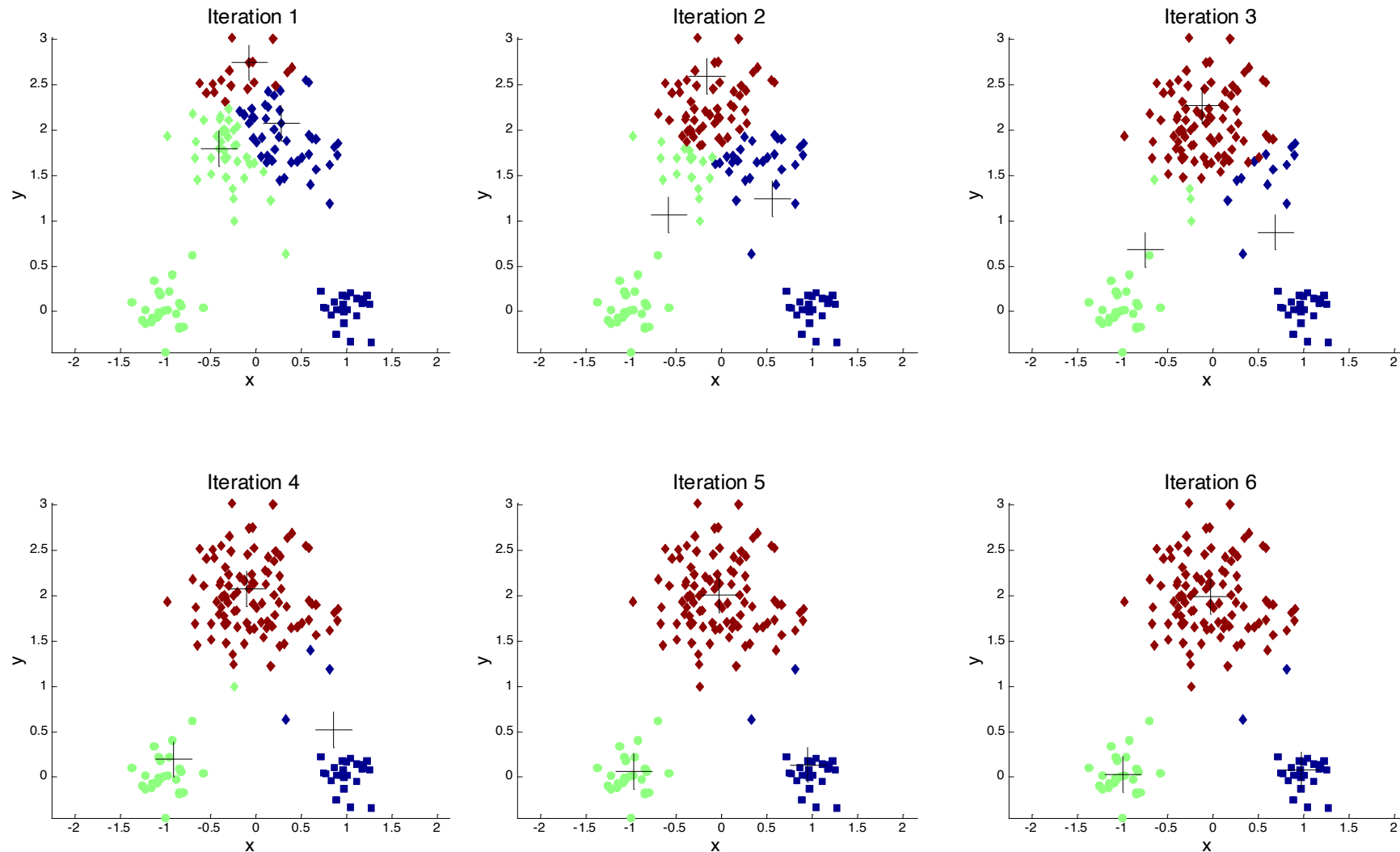
- 
- 1: Select  $K$  points as the initial centroids.
  - 2: **repeat**
  - 3:   Form  $K$  clusters by assigning all points to the closest centroid.
  - 4:   Recompute the centroid of each cluster.
  - 5: **until** The centroids don't change
-



# Example of K-means Clustering



# Example of K-means Clustering

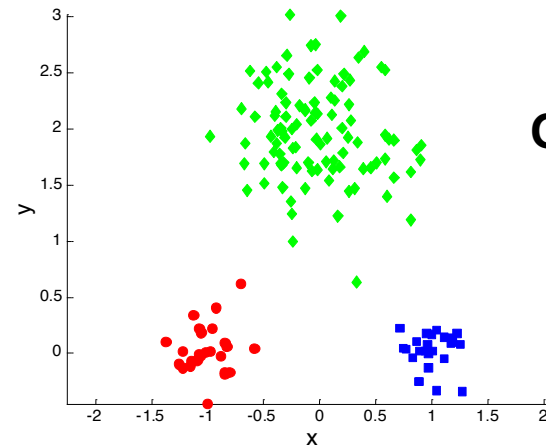


# K-means Clustering – Details

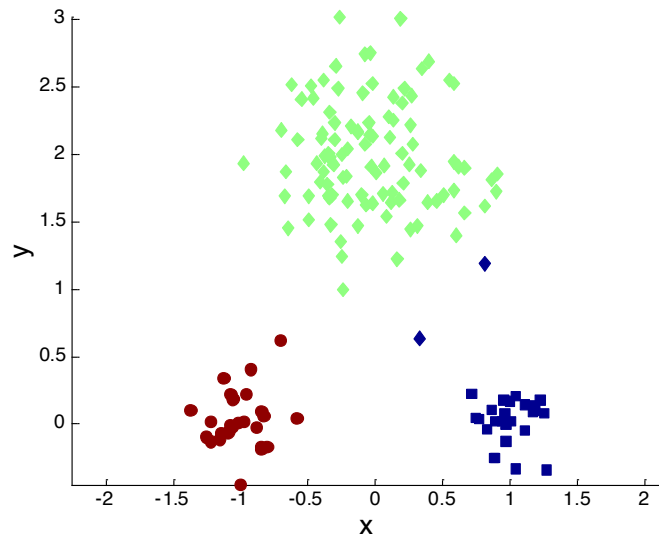
---

- Simple iterative algorithm.
  - Choose initial centroids;
  - repeat {assign each point to a nearest centroid; re-compute cluster centroids}
  - until centroids stop changing.
- Initial centroids are often chosen randomly.
  - Clusters produced can vary from one run to another
- The centroid is (typically) the mean of the points in the cluster, but other definitions are possible.
- K-means will converge for common proximity measures with appropriately defined centroid
- Most of the convergence happens in the first few iterations.
  - Often the stopping condition is changed to 'Until relatively few points change clusters'

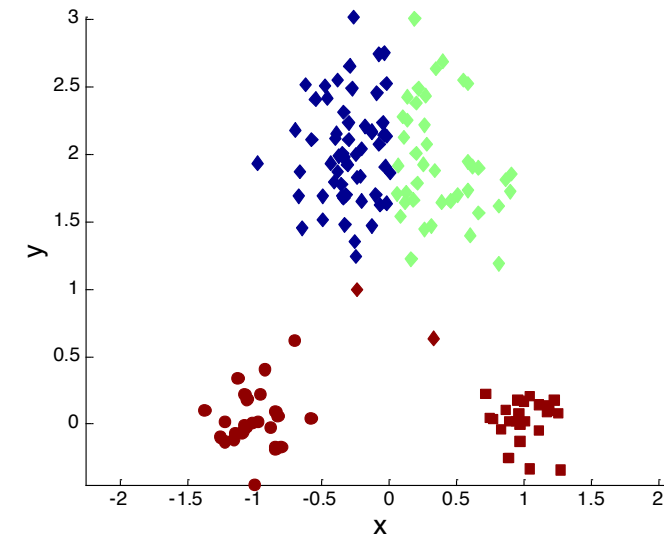
# Two different K-means Clusterings



Original Points

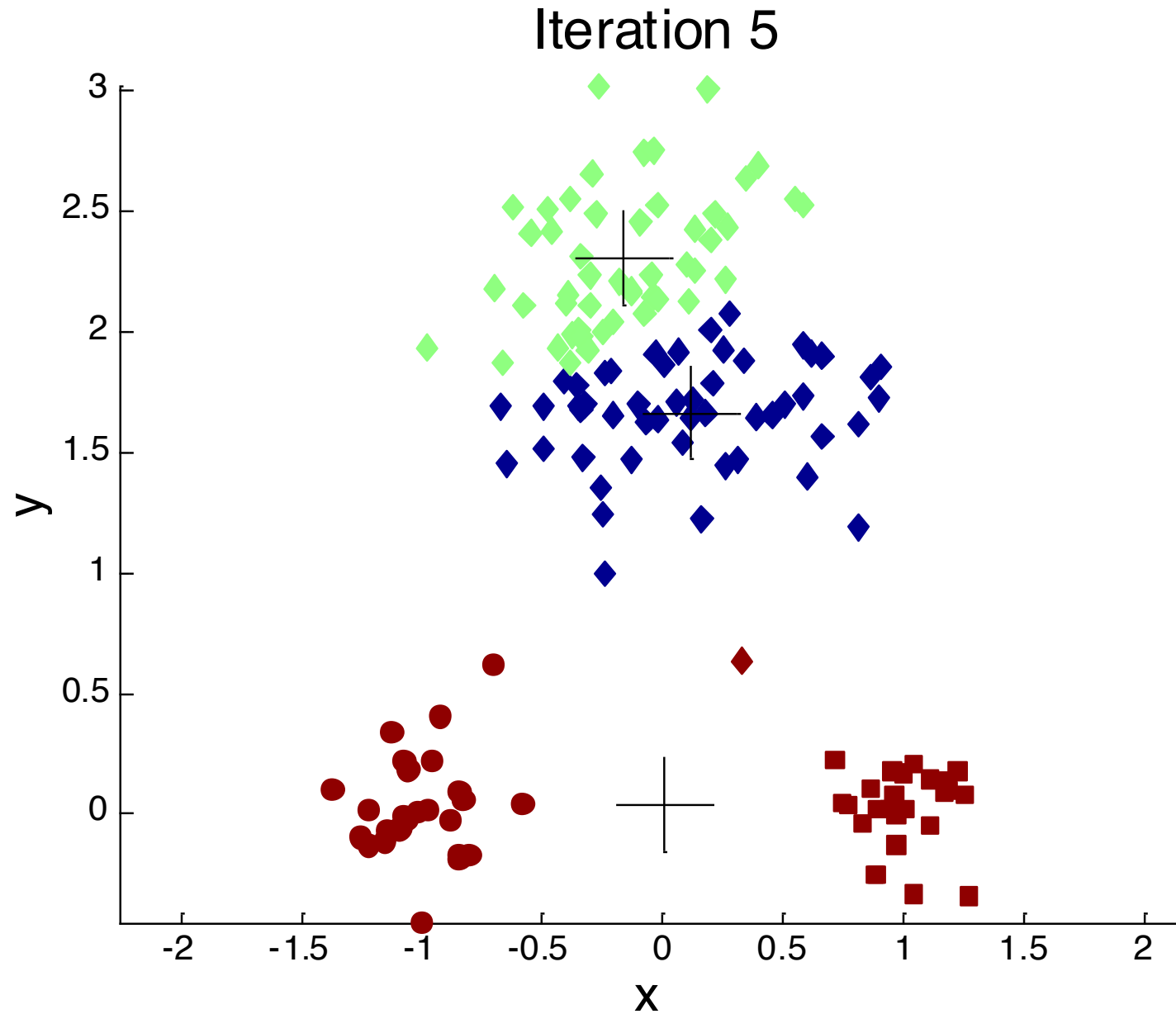


Optimal Clustering

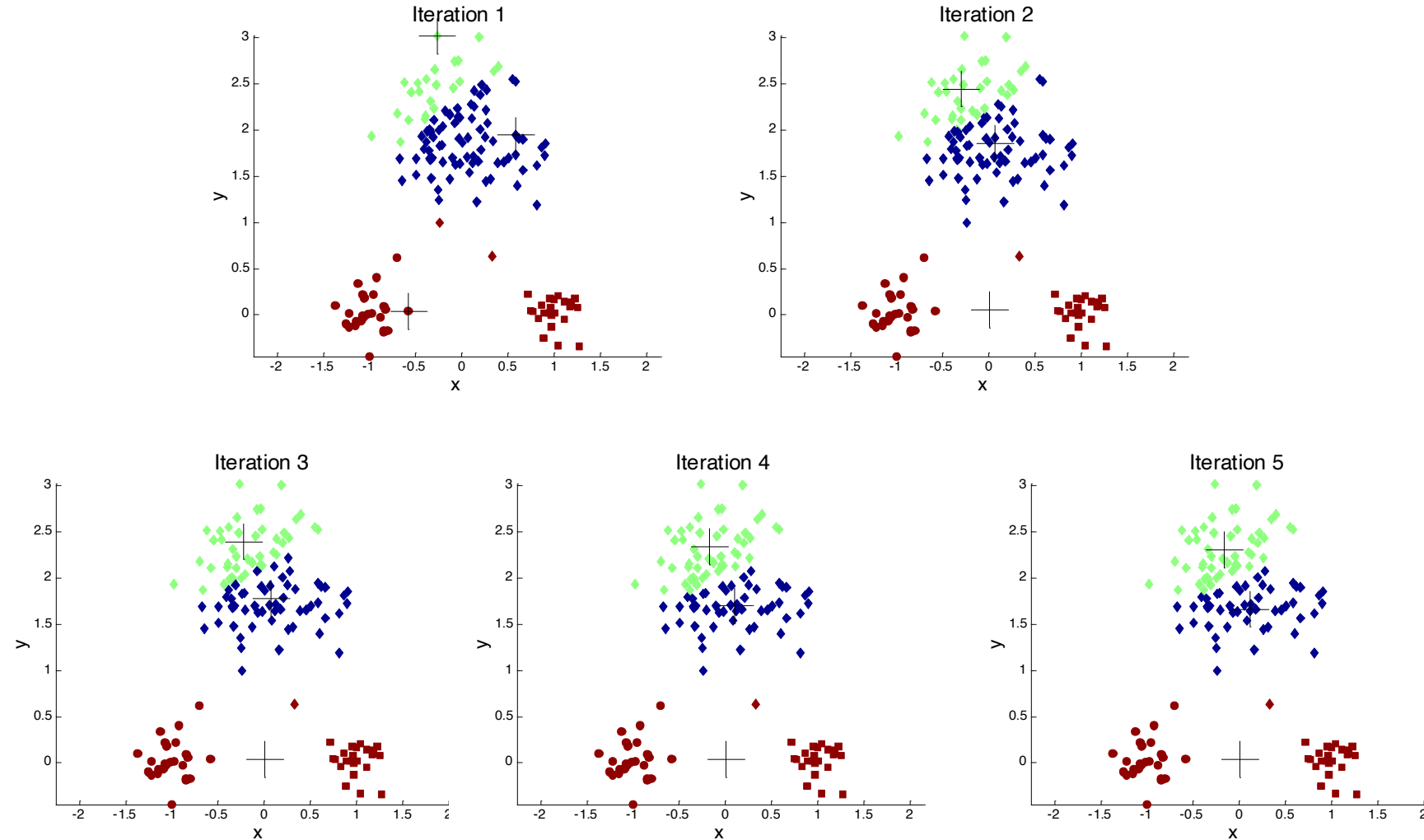


Sub-optimal Clustering

# Importance of Choosing Initial Centroids ...



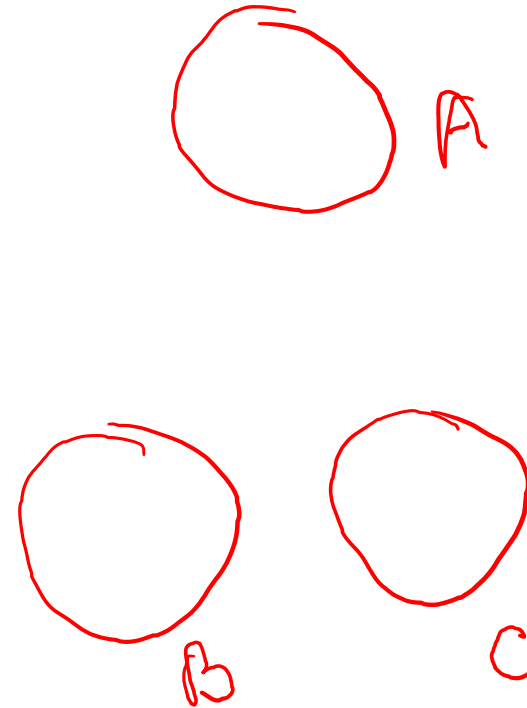
# Importance of Choosing Initial Centroids ...



# Importance of Choosing Initial Centroids

---

- Depending on the choice of initial centroids, B and C may get merged or remain separate



# Problems with Selecting Initial Points

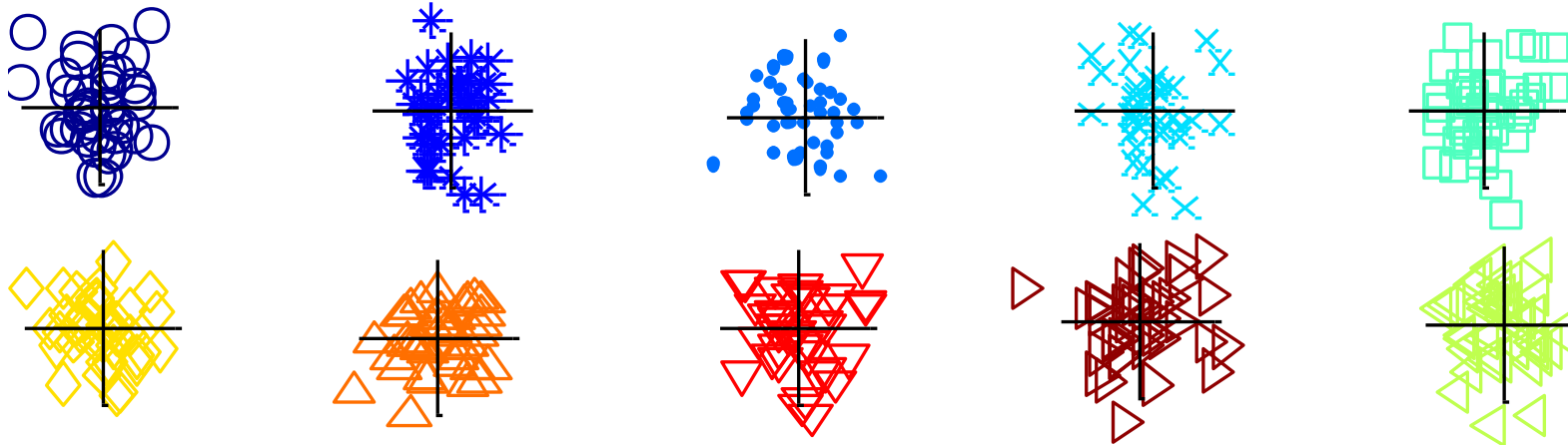
- If there are  $K$  'real' clusters then the chance of selecting one centroid from each cluster is small.
  - Chance is relatively small when  $K$  is large
  - If clusters are the same size,  $n$ , then

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

- 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

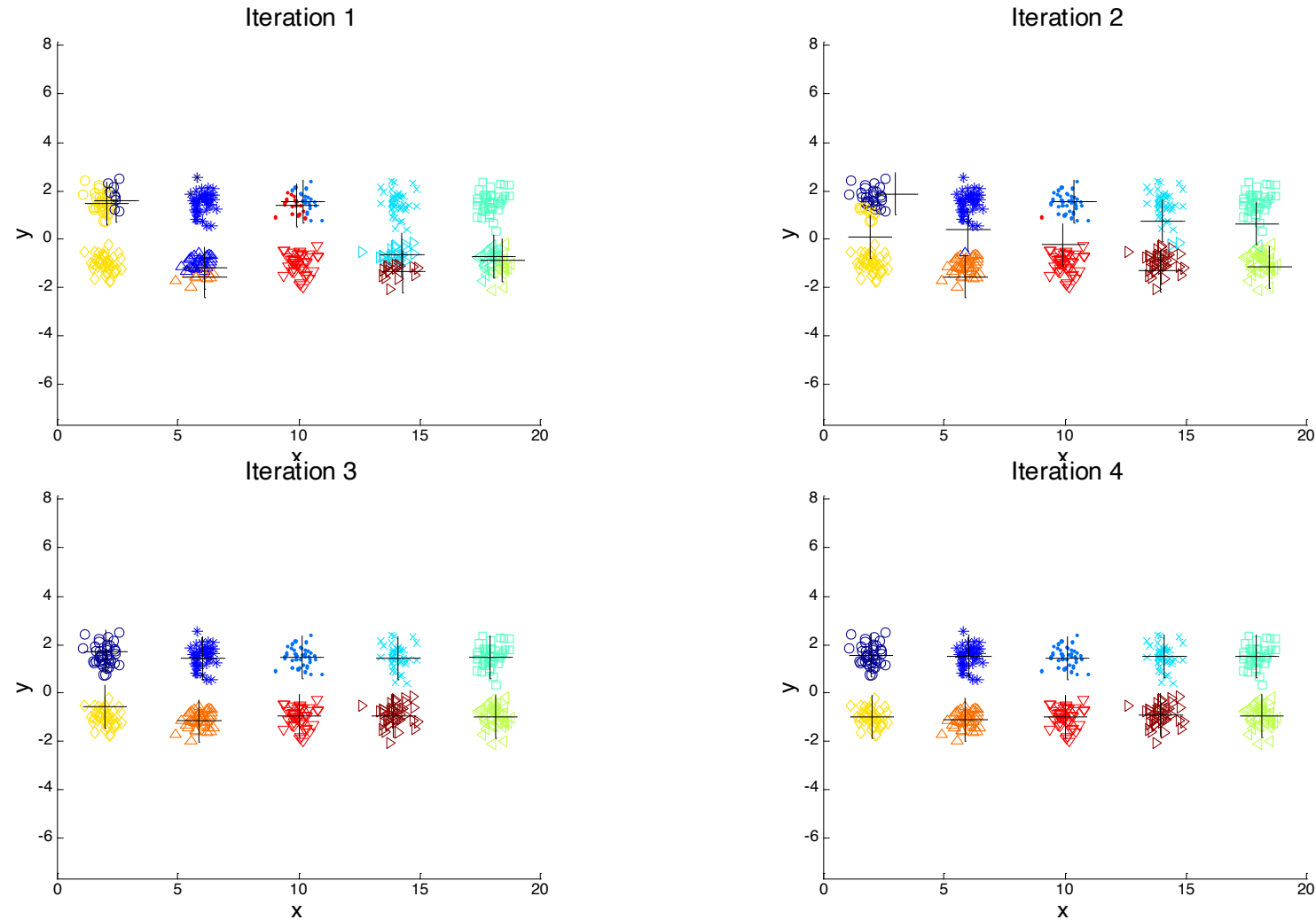


# 10 Clusters Example



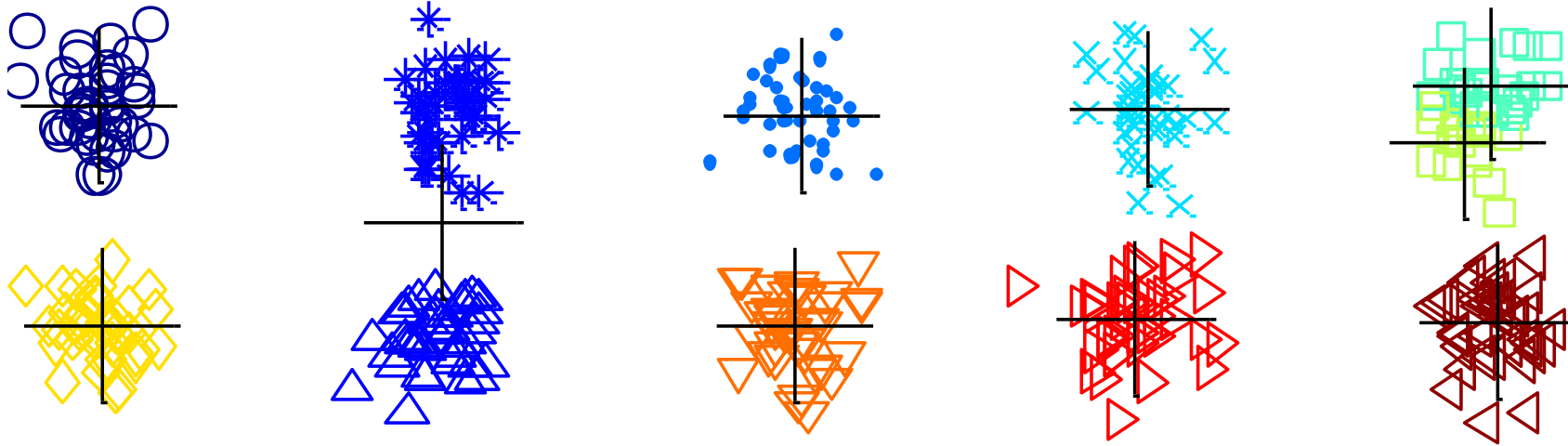
**Starting with two initial centroids in one cluster of each pair of clusters**

# 10 Clusters Example



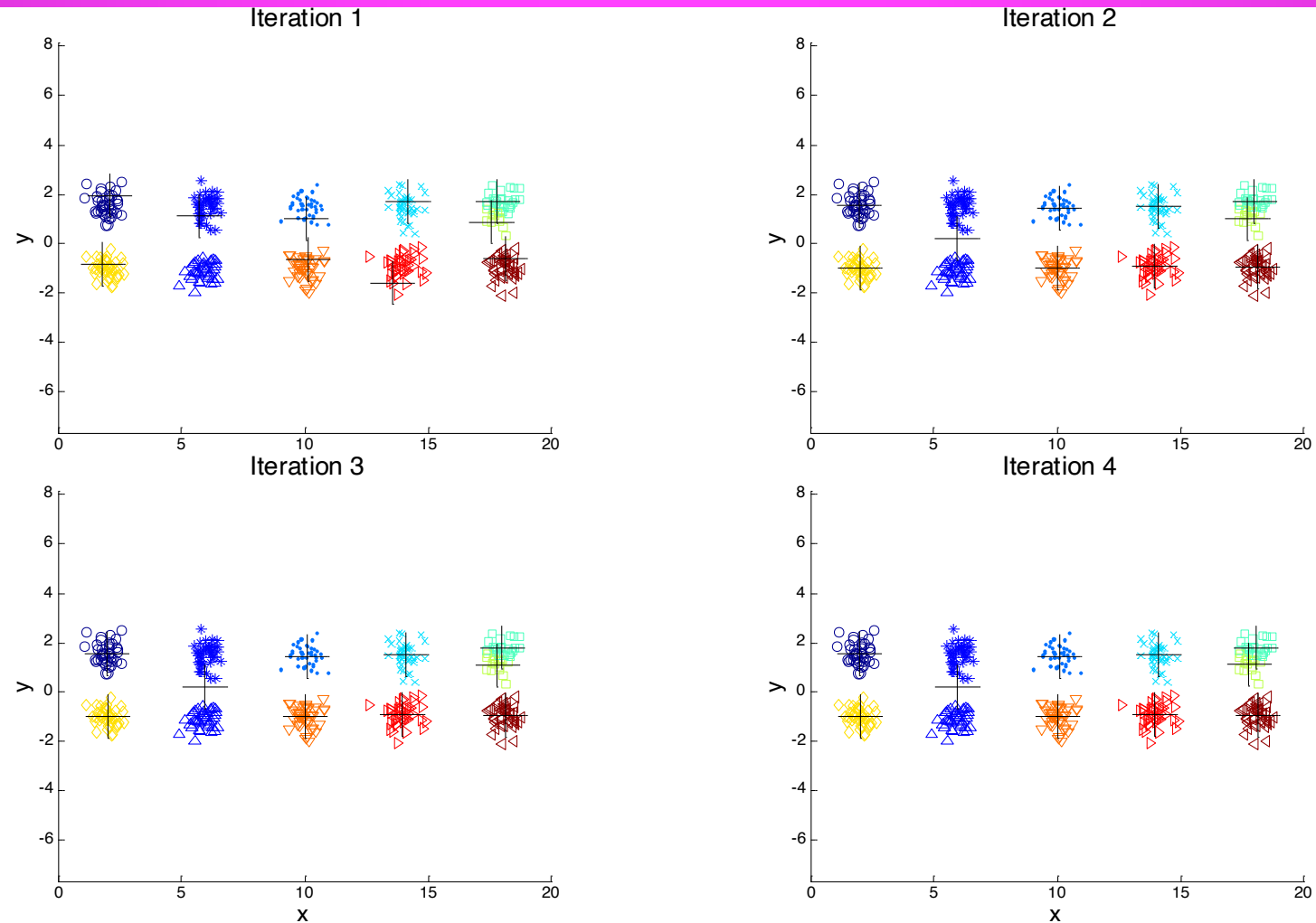
**Starting with two initial centroids in one cluster of each pair of clusters**

# 10 Clusters Example



**Starting with some pairs of clusters having three initial centroids, while other have only one.**

# 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
- Use some strategy to select the  $k$  initial centroids and then select among these initial centroids
  - Select most widely separated
    - ◆ K-means++ is a robust way of doing this selection
  - Use hierarchical clustering to determine initial centroids
- ~~Bisecting K-means~~
  - ~~Not as susceptible to initialization issues~~

# K-means++

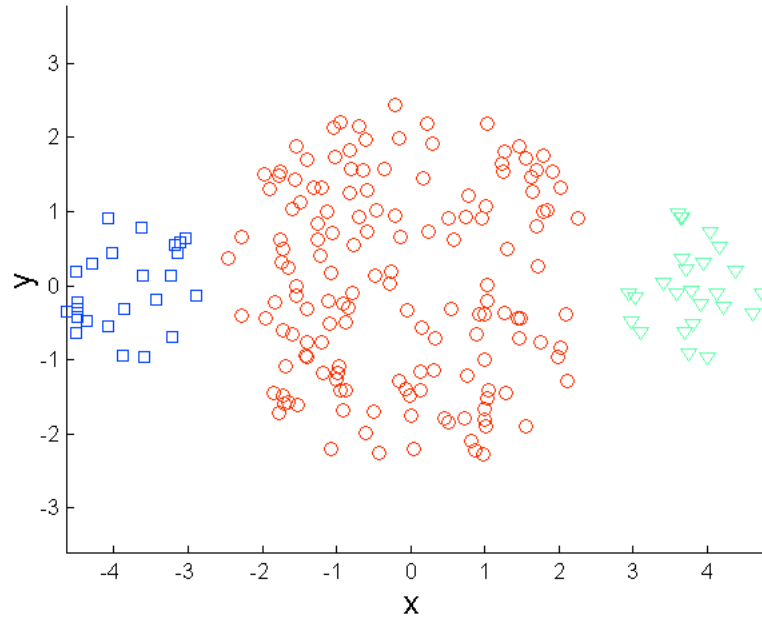
- This approach can be slower than random initialization, but very consistently produces better results in terms of SSE
  - The k-means++ algorithm guarantees an approximation ratio  $O(\log k)$  in expectation, where  $k$  is the number of centers
- To select a set of initial centroids,  $C$ , perform the following
  1. Select an initial point at random to be the first centroid
  2. For  $k - 1$  steps
    3. For each of the  $N$  points,  $x_i$ ,  $1 \leq i \leq N$ , find the minimum squared distance to the currently selected centroids,  $C_1, \dots, C_j$ ,  $1 \leq j < k$ , i.e.,  $\min_j d^2(C_j, x_i)$
    4. Randomly select a new centroid by choosing a point with probability proportional to  $\frac{\min_j d^2(C_j, x_i)}{\sum_i \min_j d^2(C_j, x_i)}$  is
  5. End For

# Limitations of K-means

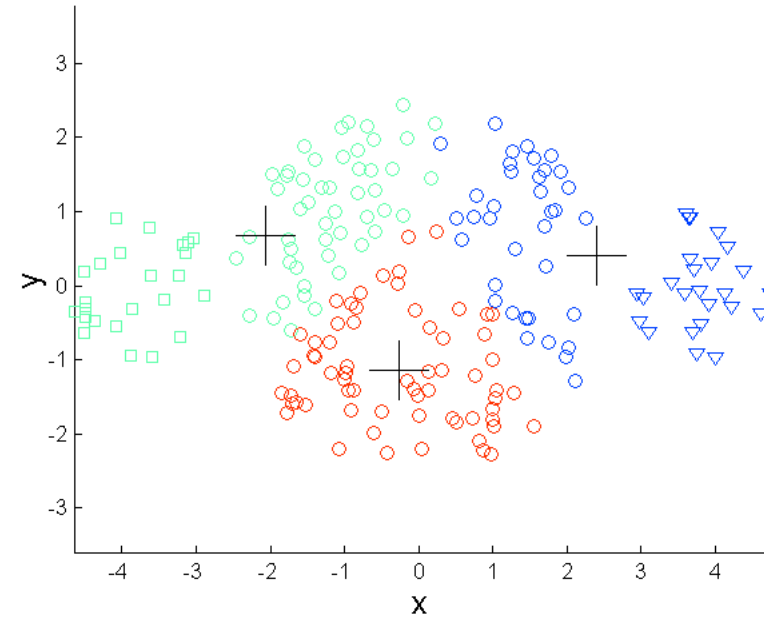
---

- K-means has problems when clusters are of differing
  - Sizes
  - Densities
  - Non-globular shapes
- K-means has problems when the data contains outliers.
  - One possible solution is to remove outliers before clustering

# Limitations of K-means: Differing Sizes



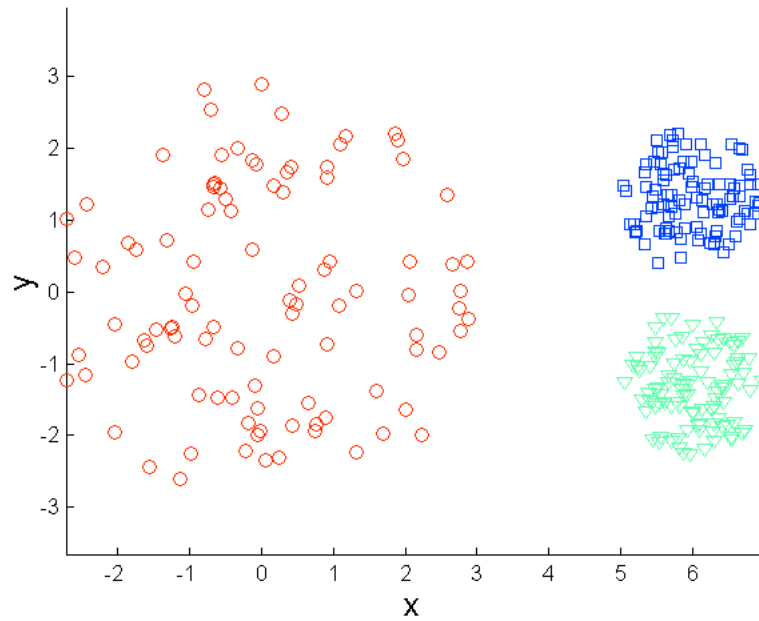
**Original Points**



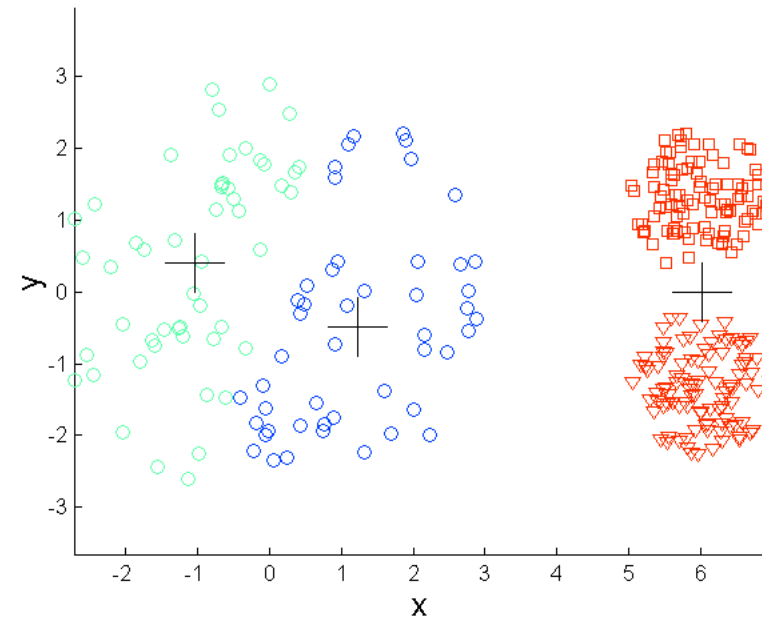
**K-means (3 Clusters)**



# Limitations of K-means: Differing Density

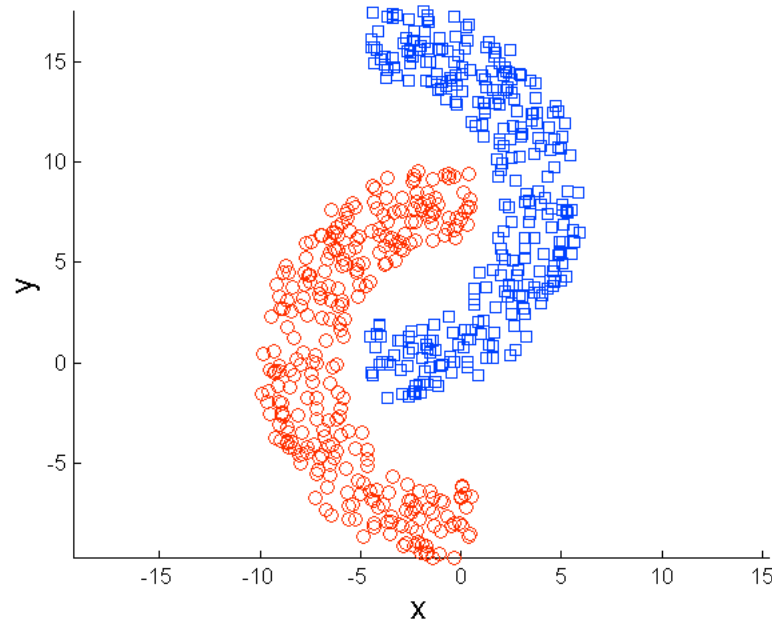


**Original Points**

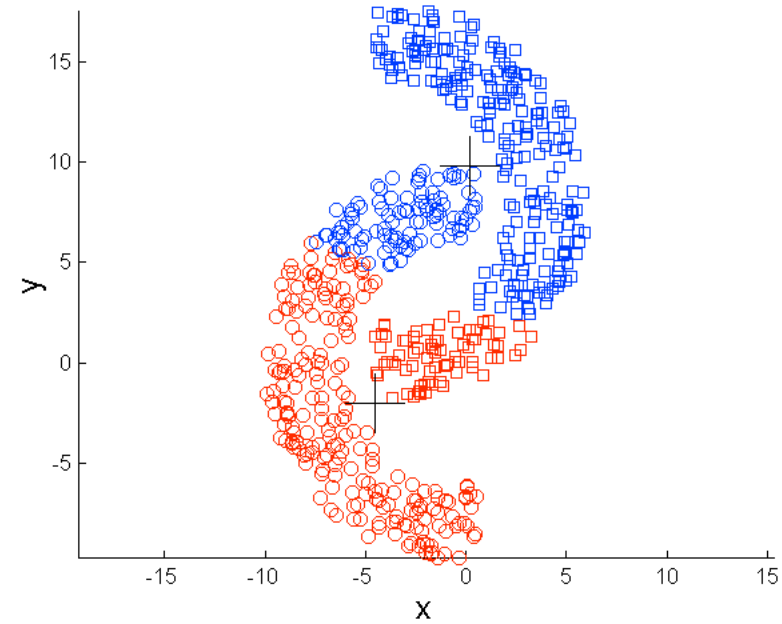


**K-means (3 Clusters)**

# Limitations of K-means: Non-globular Shapes

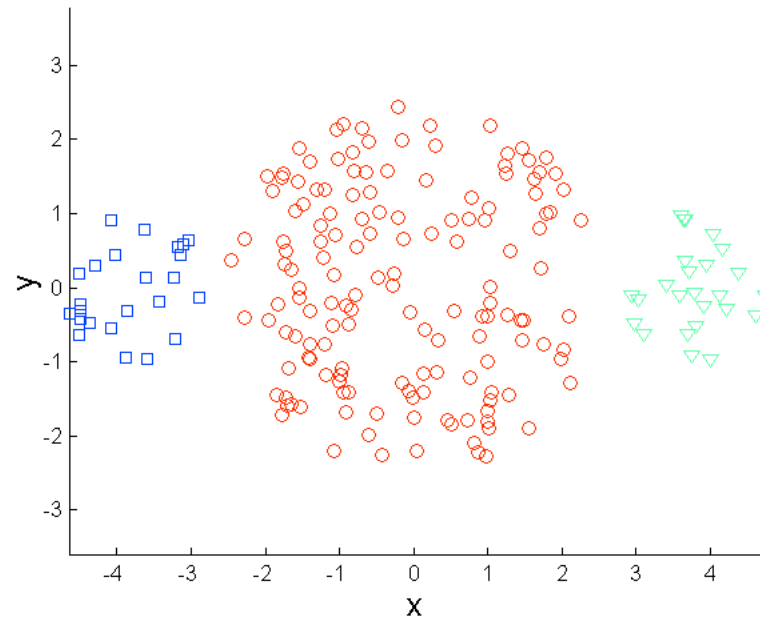


**Original Points**

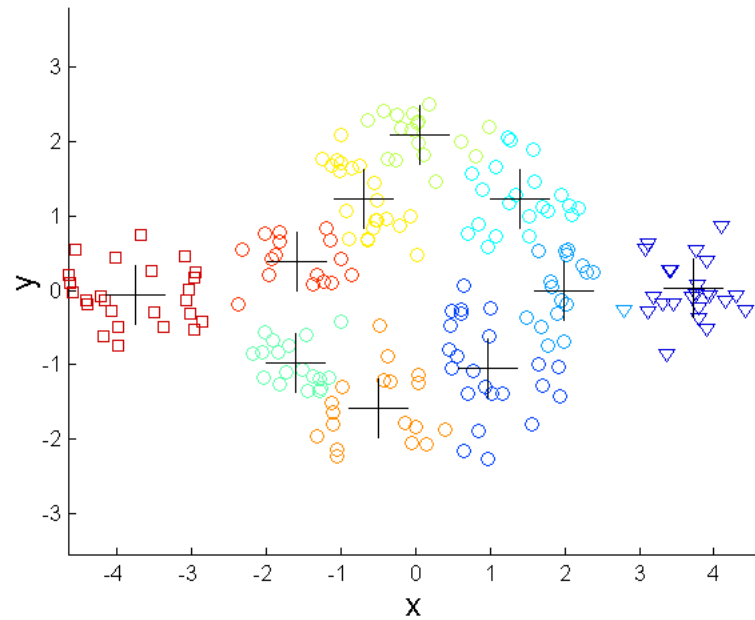


**K-means (2 Clusters)**

# Overcoming K-means Limitations



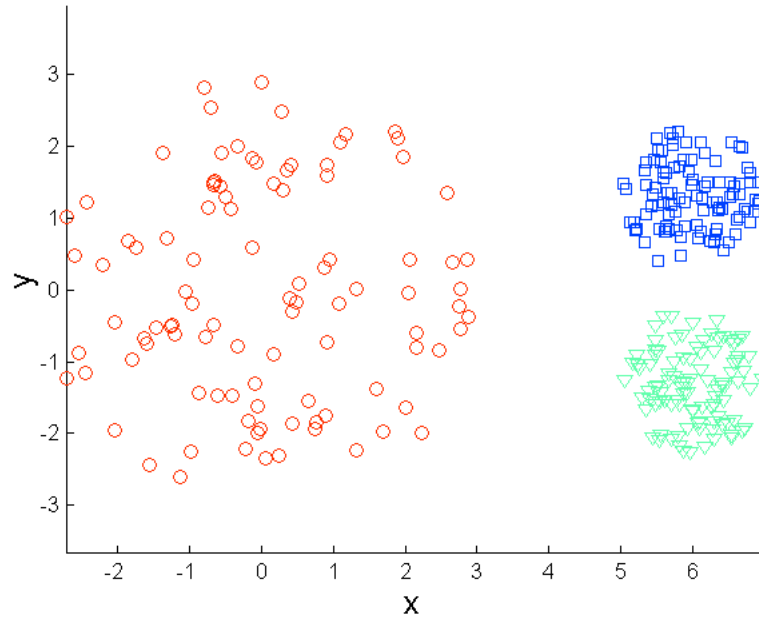
**Original Points**



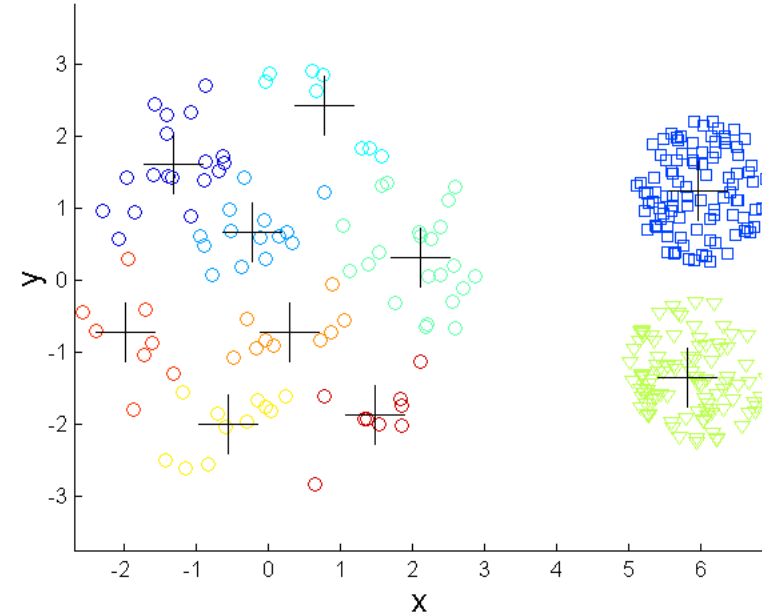
**K-means Clusters**

One solution is to find a large number of clusters such that each of them represents a part of a natural cluster. But these small clusters need to be put together in a post-processing step.

# Overcoming K-means Limitations



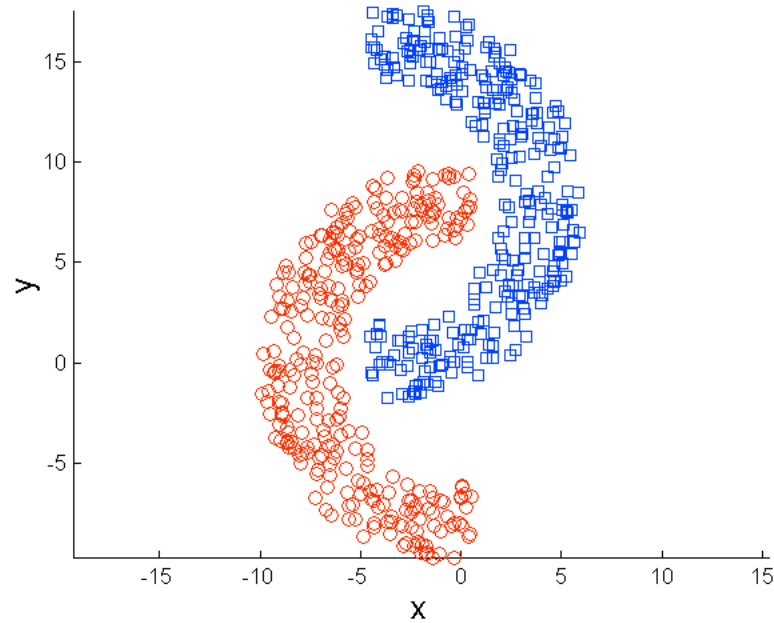
**Original Points**



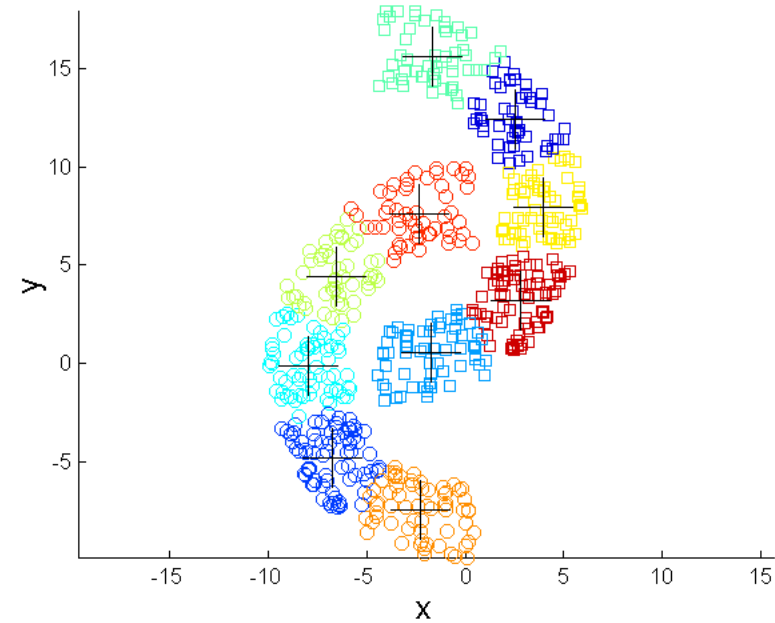
**K-means Clusters**

One solution is to find a large number of clusters such that each of them represents a part of a natural cluster. But these small clusters need to be put together in a post-processing step.

# Overcoming K-means Limitations



**Original Points**

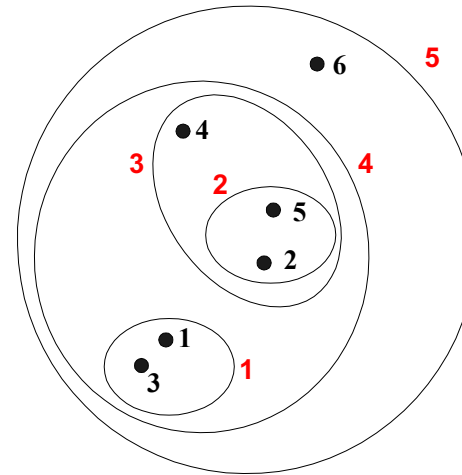
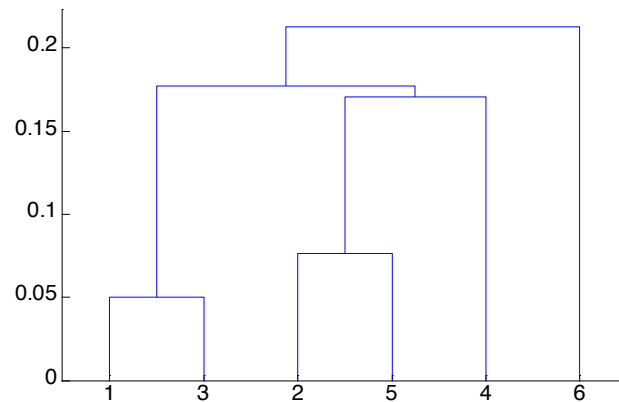


**K-means Clusters**

One solution is to find a large number of clusters such that each of them represents a part of a natural cluster. But these small clusters need to be put together in a post-processing step.

# 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 (bottom-up):
    - ◆ 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 (top-down):
    - ◆ Start with one, all-inclusive cluster
    - ◆ At each step, split a cluster until each cluster contains an individual point (or there are  $k$  clusters)
- Traditional hierarchical algorithms use a similarity or distance matrix
  - Merge or split one cluster at a time

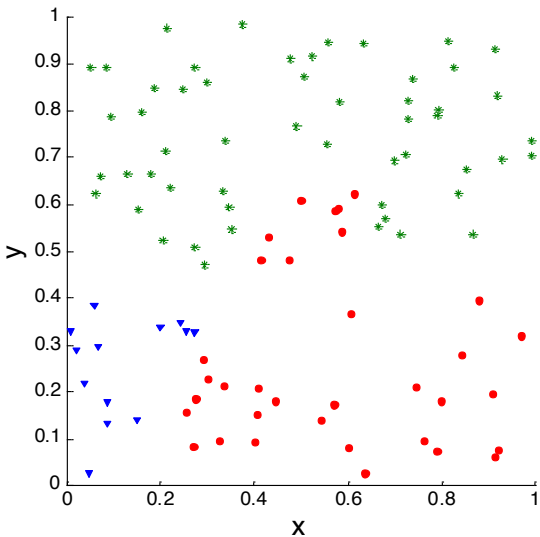
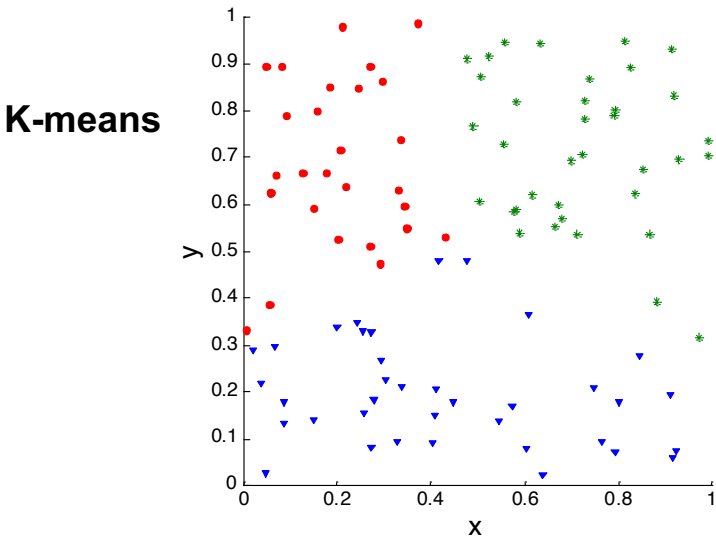
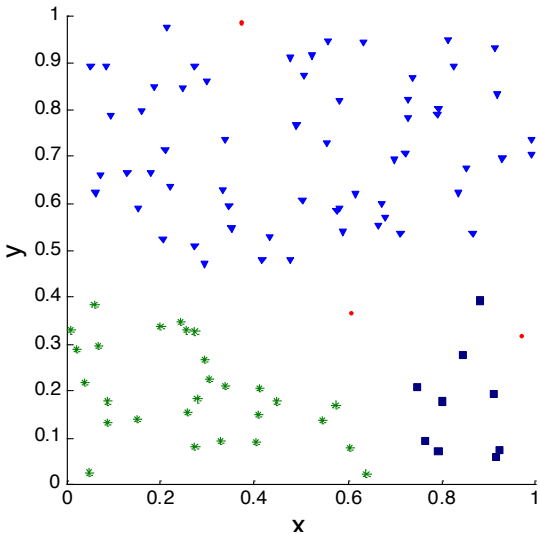
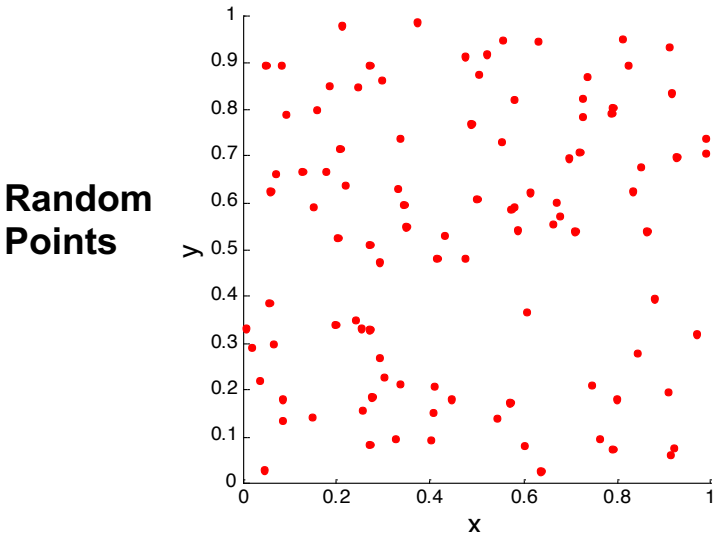


# Cluster Validity

---

- 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”!
  - In practice the clusters we find are defined by the clustering algorithm
- 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



# Measures of Cluster Validity

---

- Numerical measures that are applied to judge various aspects of cluster validity, are classified into the following two types.
  - **Supervised:** Used to measure the extent to which cluster labels match externally supplied class labels.
    - ◆ Entropy
    - ◆ Often called *external indices* because they use information external to the data
  - **Unsupervised:** Used to measure the goodness of a clustering structure *without* respect to external information.
    - ◆ Sum of Squared Error (SSE)
    - ◆ Often called *internal indices* because they only use information in the data
- You can use supervised or unsupervised measures to compare clusters or clusterings

# K-means Objective Function

- A common objective function (used with Euclidean distance measure) is Sum of Squared Error (SSE)
  - For each point, the error is the distance to the nearest cluster center
  - 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 centroid (mean) for cluster  $C_i$
- SSE improves in each iteration of K-means until it reaches a local or global minima.
- This is called the model's inertia: the mean squared distance between each instance and its closest centroid.

# Unsupervised Measures: Cohesion and Separation

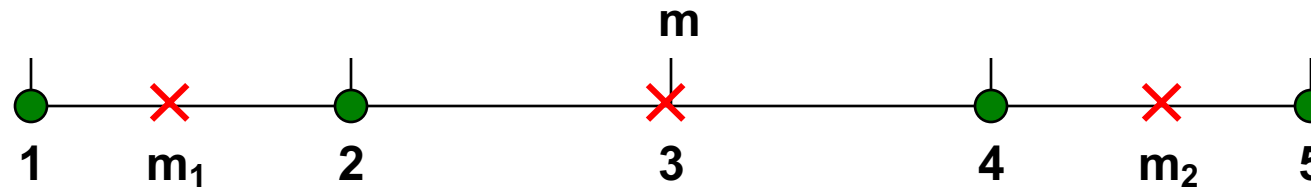
- **Cluster Cohesion:** Measures how closely related are objects in a cluster
  - Example: SSE
- **Cluster Separation:** Measure how distinct or well-separated a cluster is from other clusters
- **Example: Squared Error**
  - Cohesion is measured by the within cluster sum of squares (SSE)
$$SSE = \sum_i \sum_{x \in C_i} (x - m_i)^2$$
  - Separation is measured by the between cluster sum of squares
$$SSB = \sum_i |C_i| (m - m_i)^2$$

Where  $|C_i|$  is the size of cluster  $i$ ,  $m$  is the global average point,  $m_i$  is the centroid of cluster  $i$ , and  $x$  is a point in the cluster

# Unsupervised Measures: Cohesion and Separation

- Example: SSE

- $SSB + SSE = \text{constant}$



**K=1 cluster:**  $SSE = (1 - 3)^2 + (2 - 3)^2 + (4 - 3)^2 + (5 - 3)^2 = 10$

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

$$Total = 10 + 0 = 10$$

**K=2 clusters:**  $SSE = (1 - 1.5)^2 + (2 - 1.5)^2 + (4 - 4.5)^2 + (5 - 4.5)^2 = 1$

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

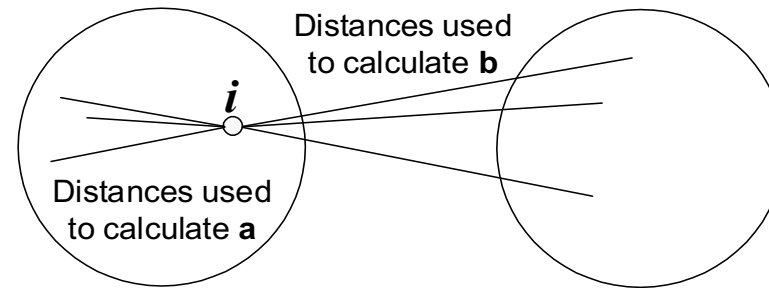
$$Total = 1 + 9 = 10$$

# Unsupervised 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

$$s = (b - a) / \max(a, b)$$

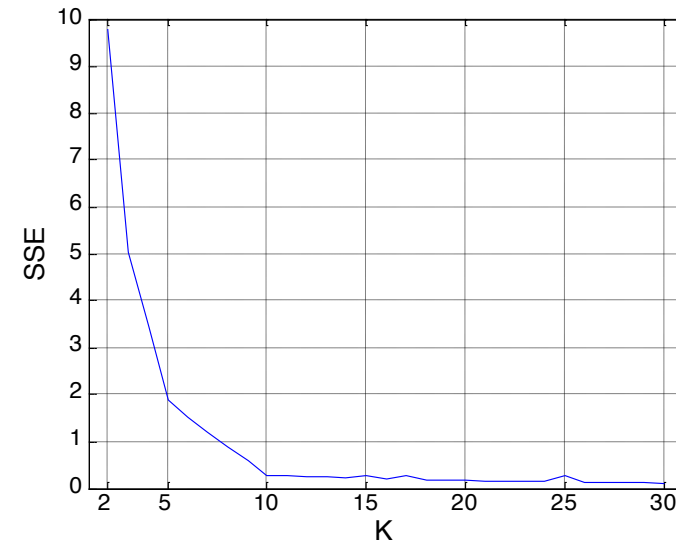
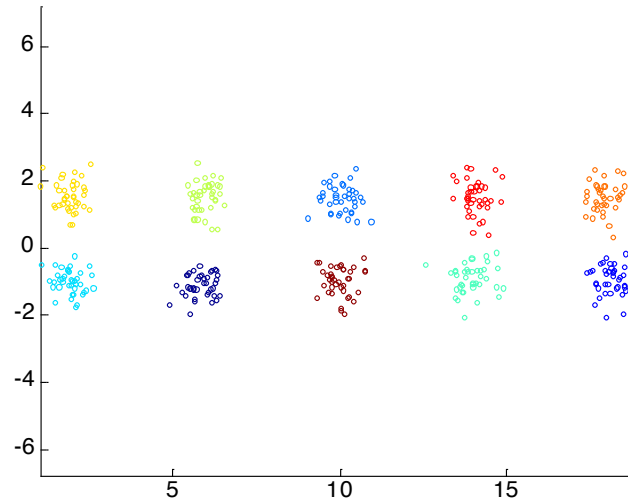
- Value can vary between -1 and 1
- Typically ranges between 0 and 1.
- The closer to 1 the better.



- Can calculate the average silhouette coefficient for a cluster or a clustering

# Determining the Correct Number of Clusters

- SSE is good for comparing two clusterings or two clusters
- SSE can also be used to estimate the number of clusters





# Selecting the Number of Clusters

Can use the elbow method of inertia, similar to elbow method of PCA.

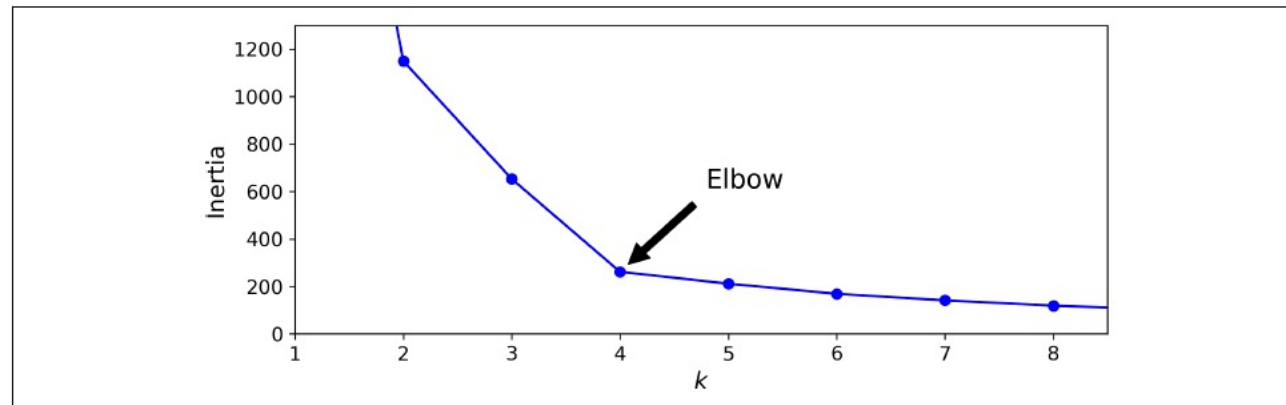
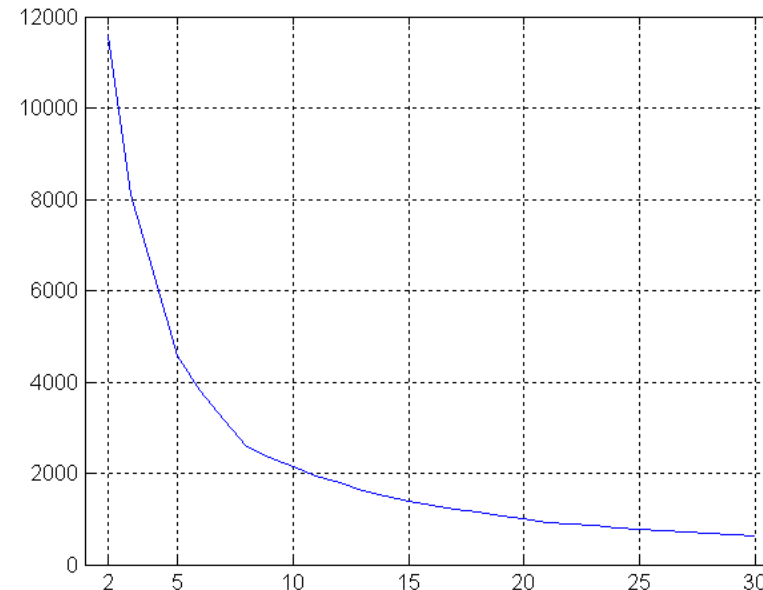
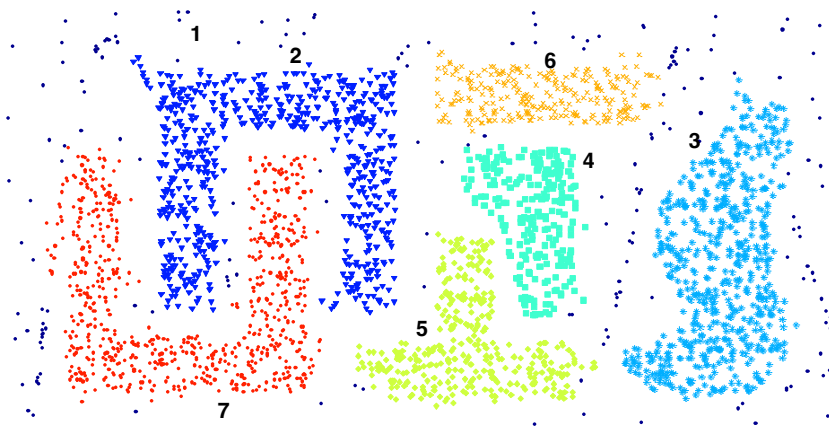


Figure 9-8. Selecting the number of clusters  $k$  using the “elbow rule”

# Determining the Correct Number of Clusters

- SSE curve for a more complicated data set



**SSE of clusters found using K-means**