

# Data Mining

---

## UNIT- IV

### Association Rule Mining

**Dr. Anshul Sharma**

Assistant Professor

Department of Computer Science & Engineering

LNM Institute of Information Technology, Jaipur, Rajasthan

# Association Analysis: Advanced Concepts

---

## Infrequent Patterns

# Infrequent Patterns

---

- An infrequent pattern is an itemset or a rule whose support is less than a *minsup* threshold.
- Example:
  - {DVDs, VCRs}
  - {Fire = yes}, {Fire = yes, Alarm=on}
- Mining infrequent patterns is challenging:
  - (1) how to identify interesting infrequent patterns.
  - (2) how to efficiently discover them in large data sets.

# Negative Patterns

---

- **Negative Itemset:** A negative itemset  $X$  is an itemset that has the following properties:

$$(1) X = A \cup \bar{B}, |\bar{B}| \geq 1$$

$$(2) s(X) \geq \textit{minsup}.$$

- **Negative Association Rule:** A negative association rule is an association rule that has the following properties:

- the rule is extracted from a negative itemset,
- the support of the rule is greater than or equal to *minsup*,
- the confidence of the rule is greater than or equal to *minconf*.

---

*Example:*

$$tea = \overline{coffee}$$

Which may suggest that people who drink tea tend to not drink coffee.

# Negatively Correlated Patterns

Let  $X = (x_1, x_2, \dots, x_k)$  denote a  $k$ -itemset and

$P(X)$  denote the probability that a transaction contains  $X$ . In association analysis, the probability is often estimated using the itemset support,  $s(X)$

- An itemset  $X$  is negatively correlated if

$$s(X) < \prod_{j=1}^k s(x_j) = s(x_1) \times s(x_2) \times \dots \times s(x_k),$$

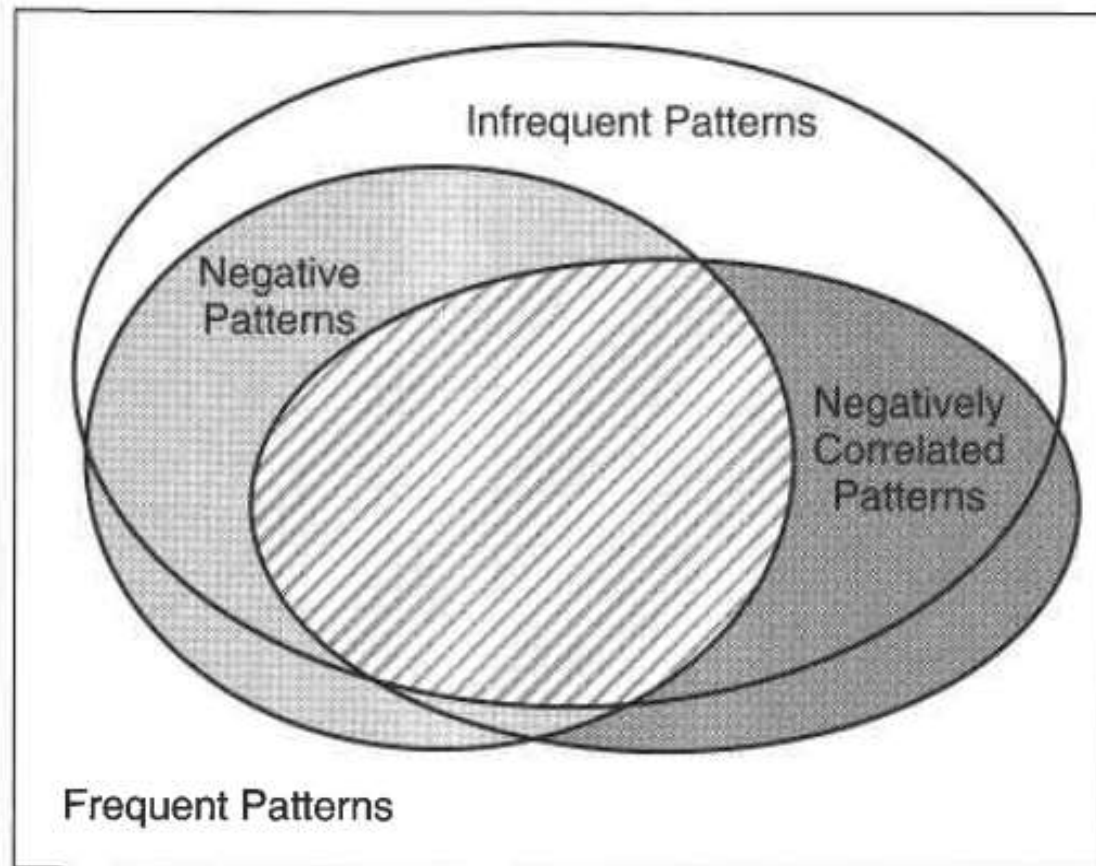
- An association rule  $X \rightarrow Y$  is negatively correlated if

$$s(X \cup Y) < s(X)s(Y),$$

- A full condition for negative correlation:

$$s(X \cup Y) < \prod_i s(x_i) \prod_j s(y_j),$$

- Comparisons among infrequent patterns, negative patterns, and negatively correlated patterns.





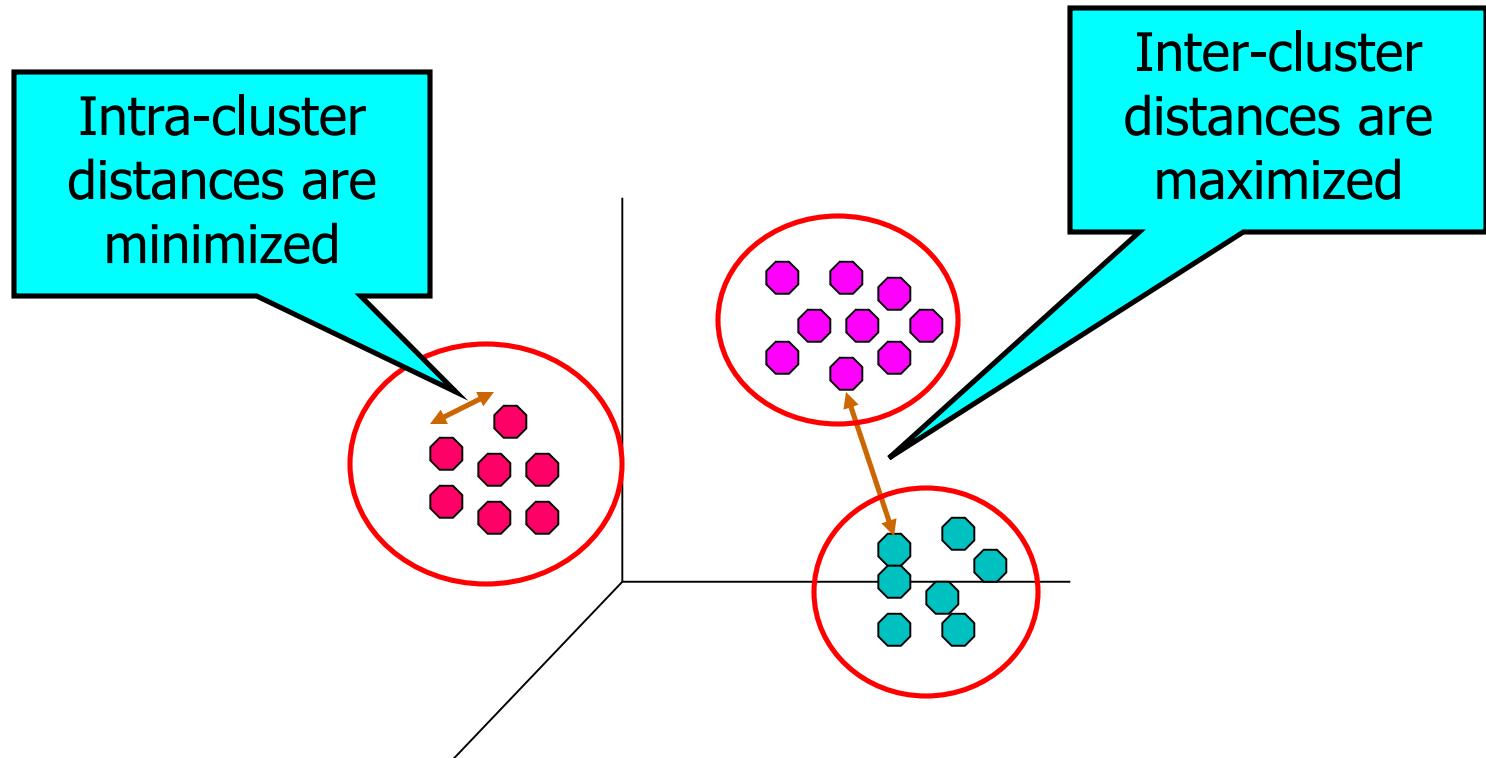
# **UNIT- V**

## **Cluster Analysis**



# 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

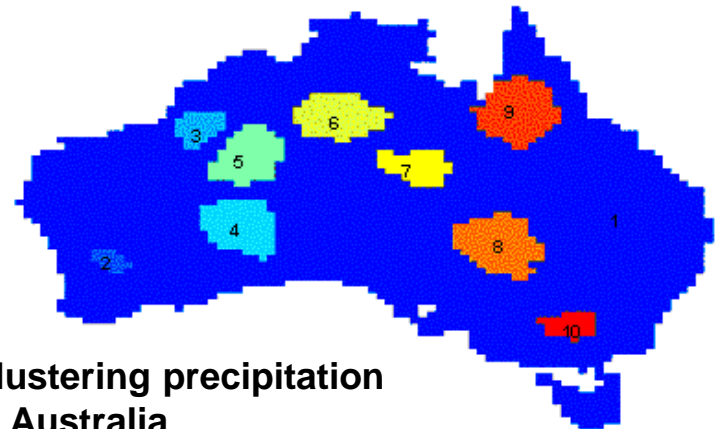
## ● Understanding

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

	<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

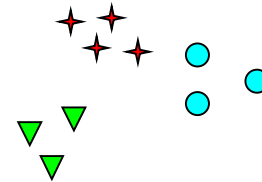


Clustering precipitation  
in Australia

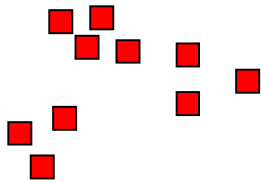
# Notion of a Cluster can be Ambiguous



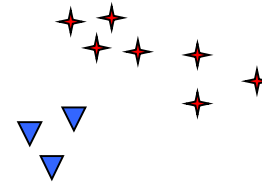
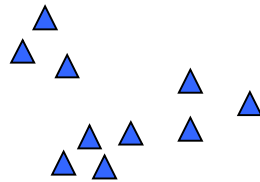
How many clusters?



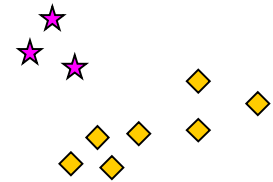
Six Clusters



Two Clusters



Four Clusters



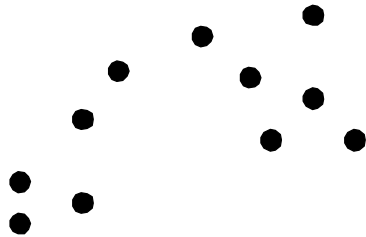
# Types of Clusterings

---

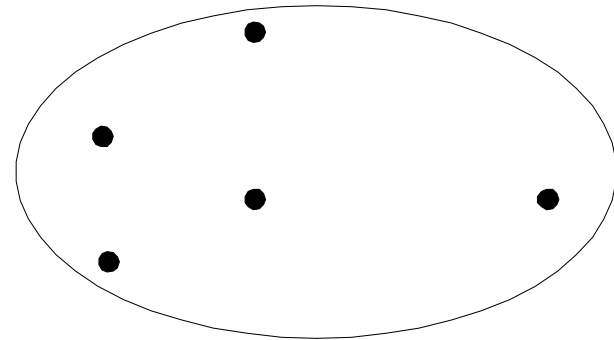
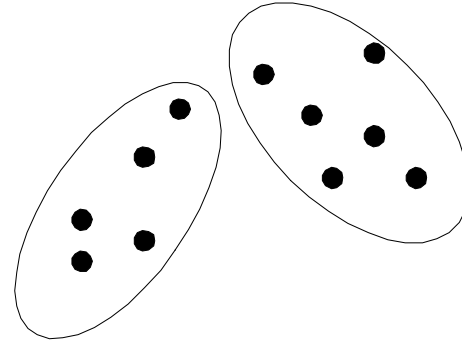
- A **clustering** is a set of clusters
- Important distinction between **hierarchical** and **partitional** sets of clusters
  - Partitional Clustering
    - ◆ A division of data objects into non-overlapping subsets (clusters)
  - Hierarchical clustering
    - ◆ A set of nested clusters organized as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

# Partitional Clustering

---

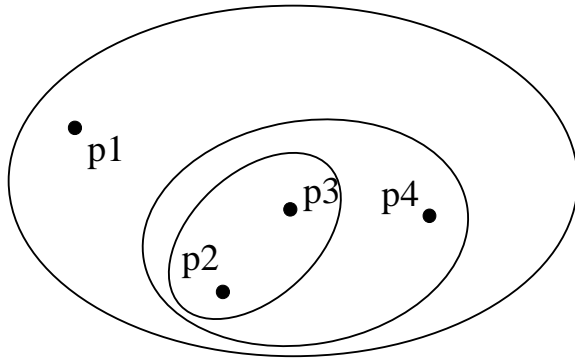


Original Points

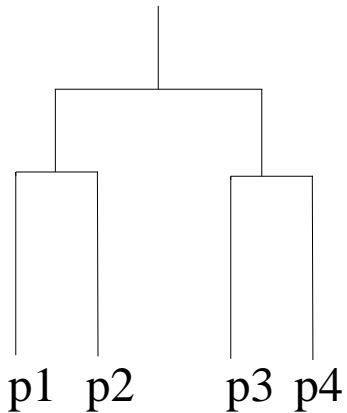
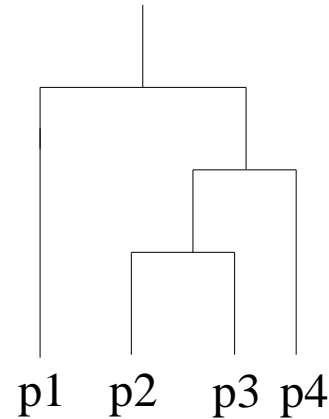


A Partitional Clustering

# Hierarchical Clustering



**Hierarchical Clustering**



**Dendrogram**

# 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.

# 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

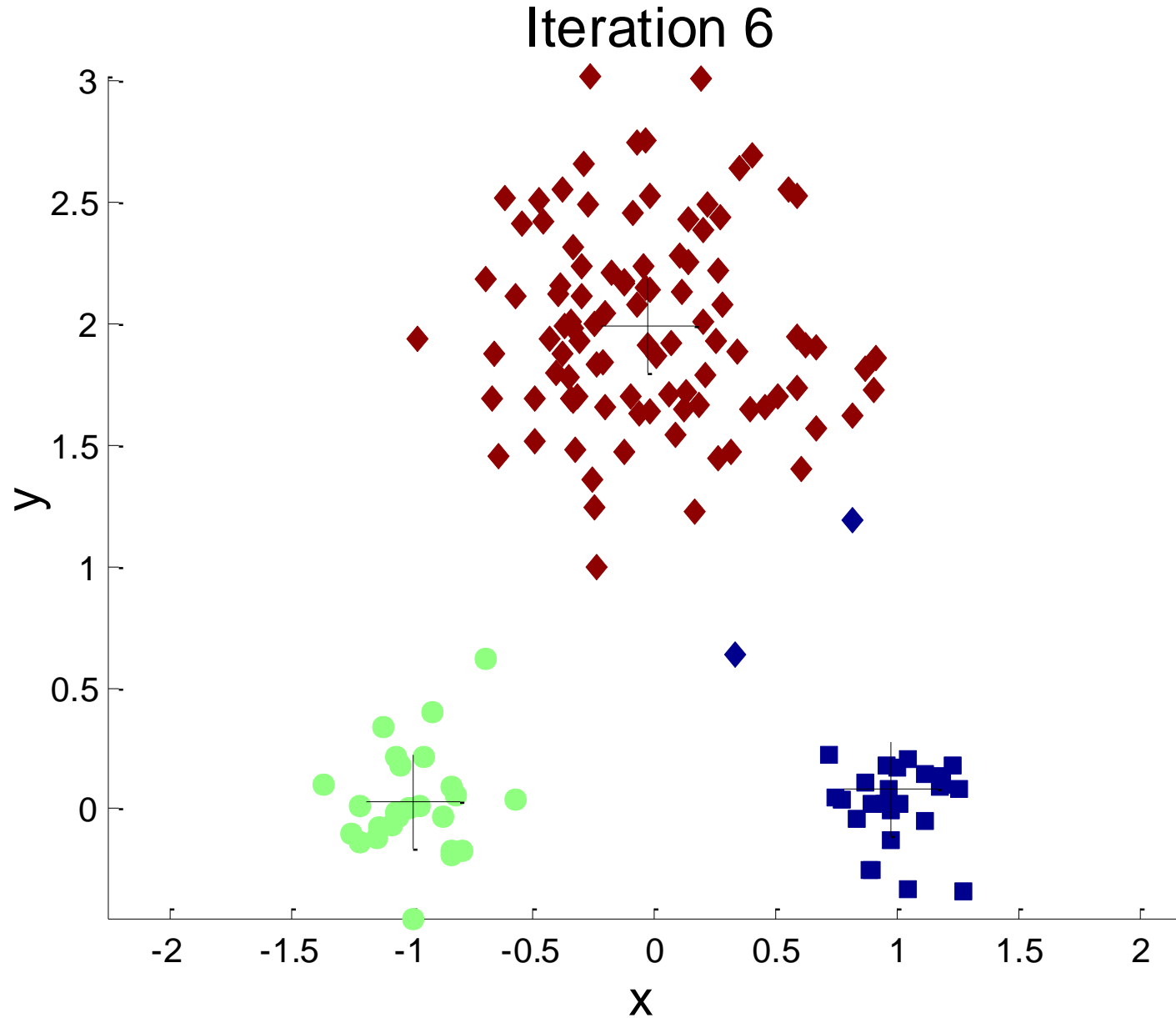
# K-means Clustering

---

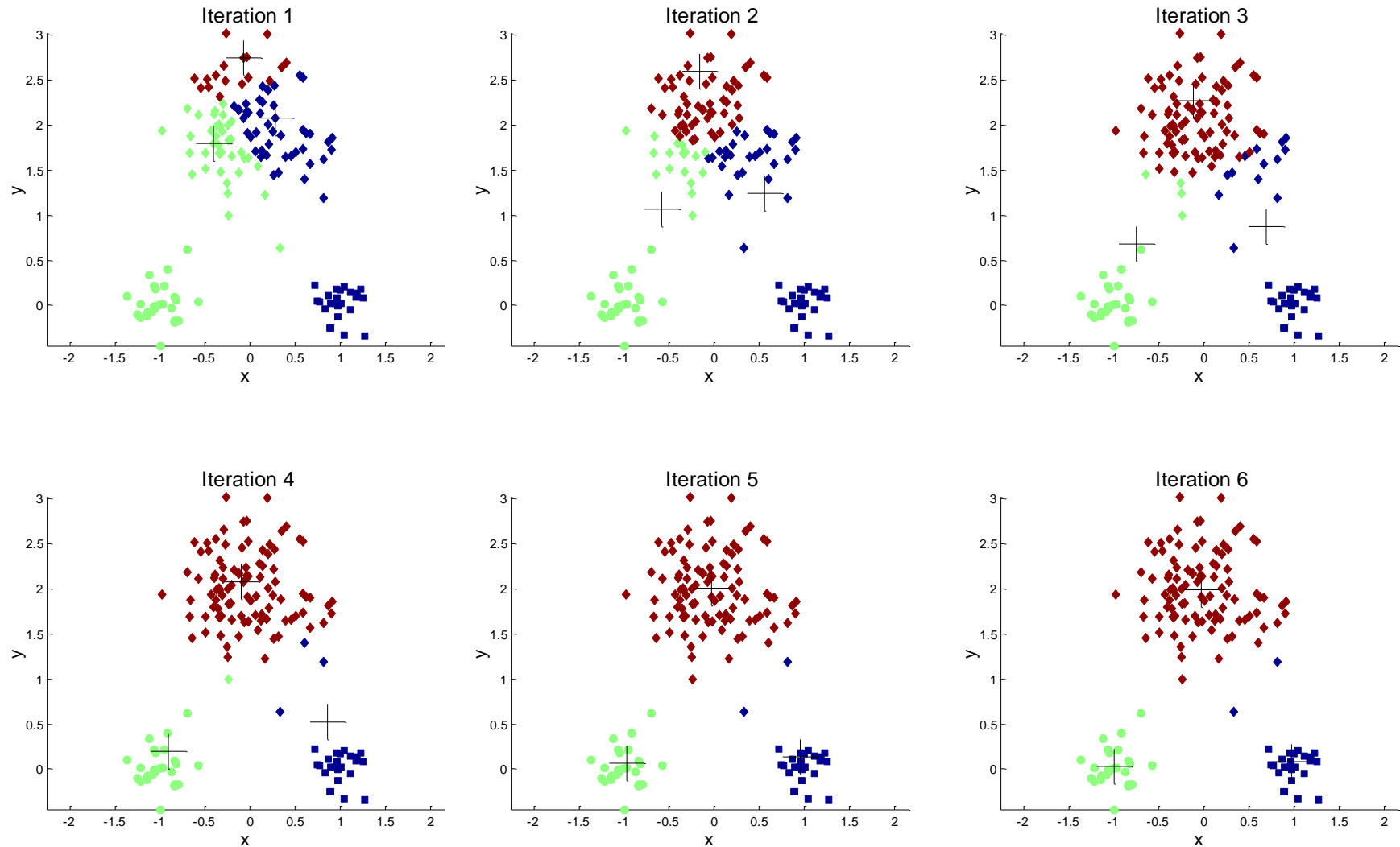
- Partitional clustering approach
- Number of clusters,  $K$ , must be specified
- Each cluster is associated with a **centroid** (center point)
- Each point is assigned to the cluster with the closest centroid
- The basic algorithm is very simple

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

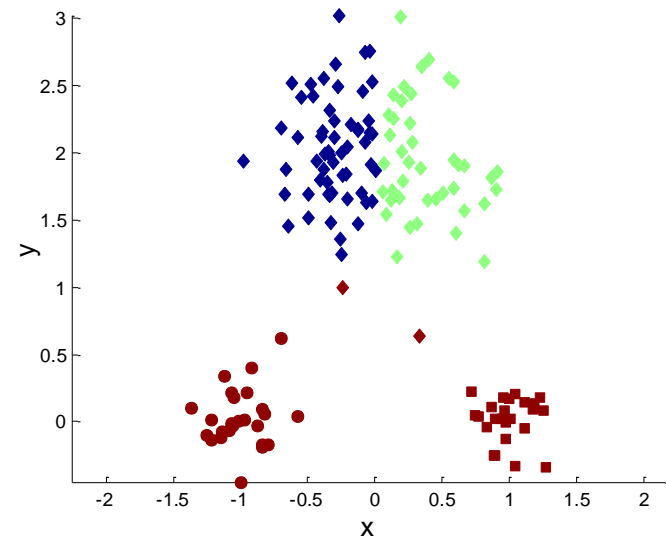
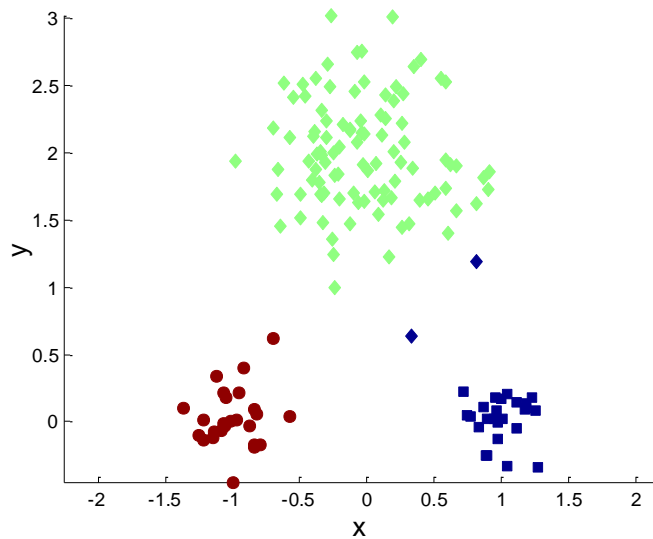
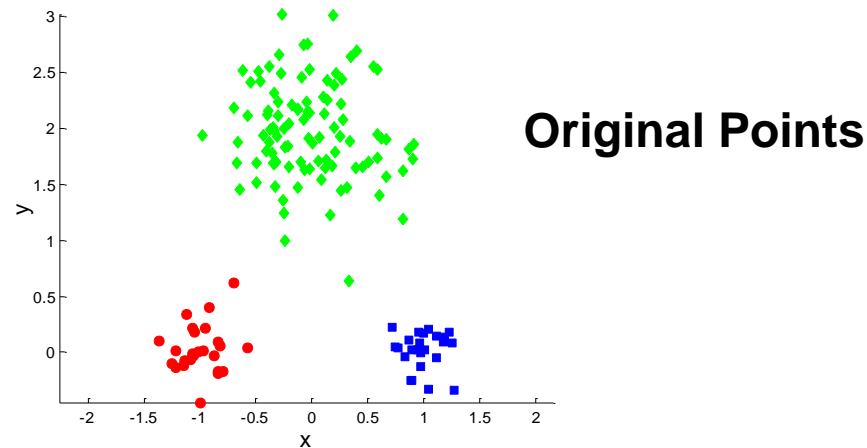
# K-means Objective Function

- 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.

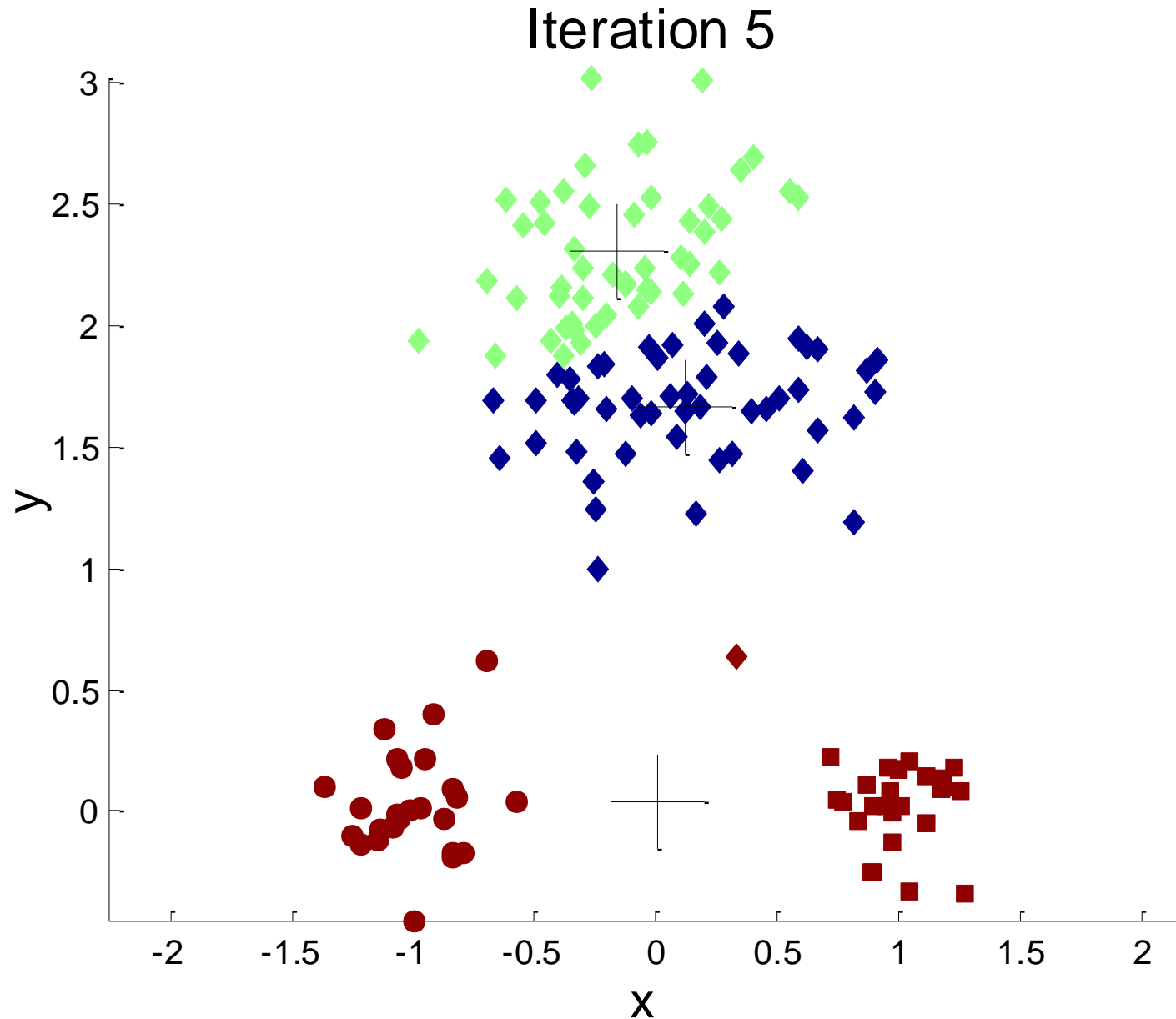
$$SSE = \sum_{i=1}^K \sum_{x \in C_i} dist^2(m_i, x)$$

- $x$  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.

# Two different K-means Clusterings

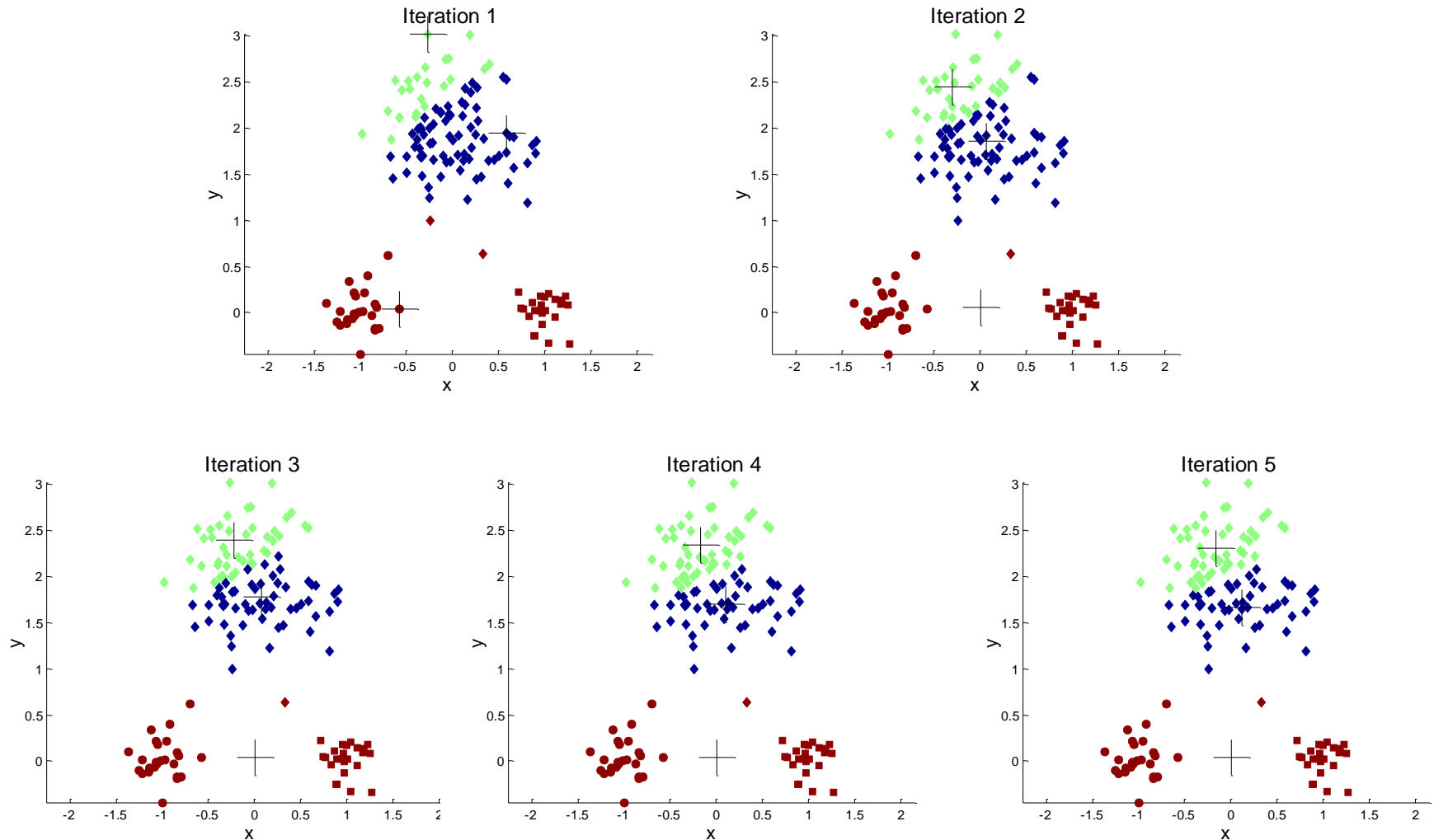


# Importance of Choosing Initial Centroids ...





# Importance of Choosing Initial Centroids ...

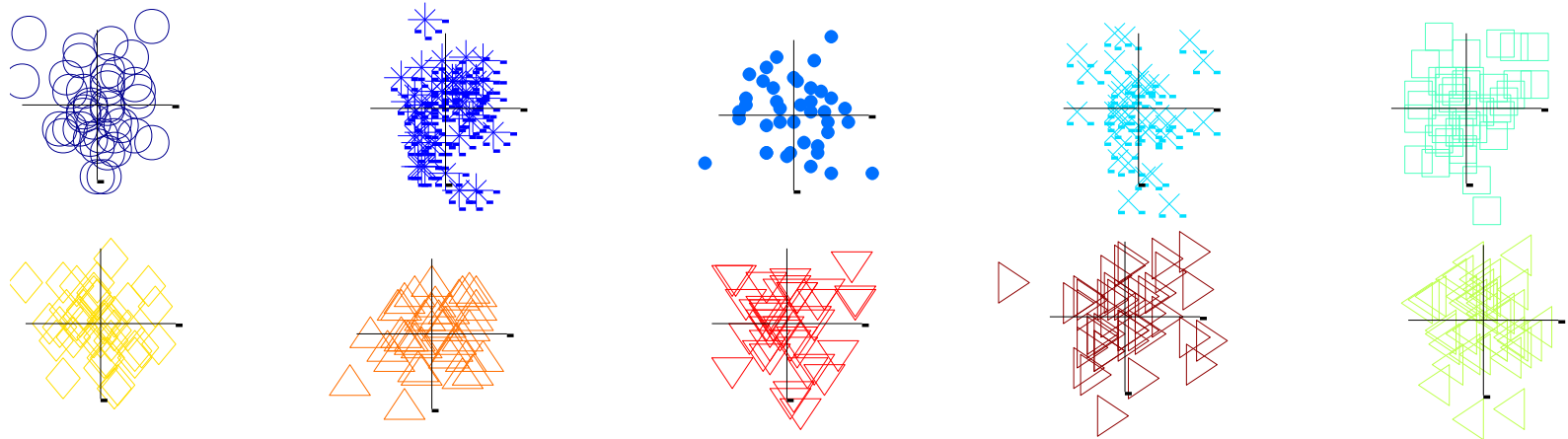


# Importance of Choosing Initial Centroids

---

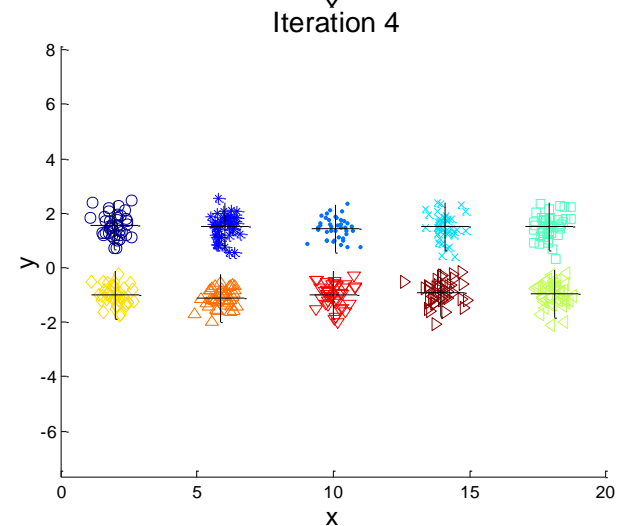
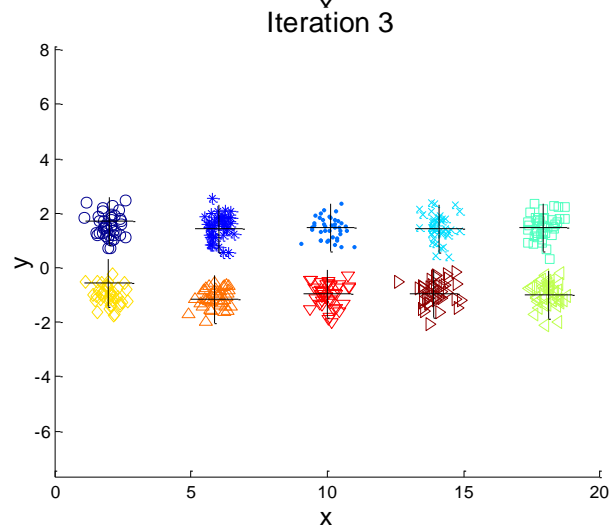
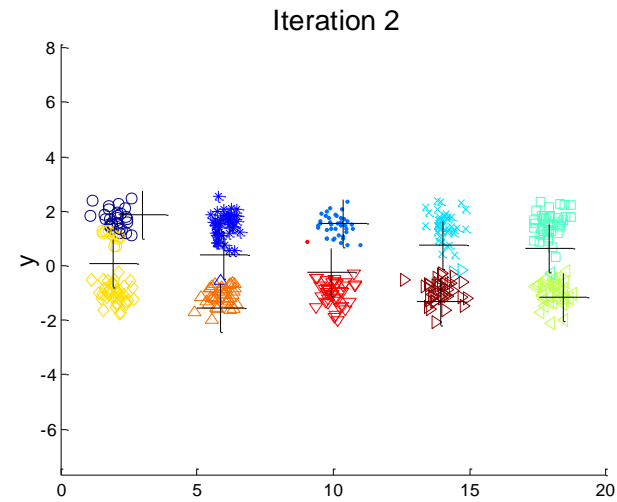
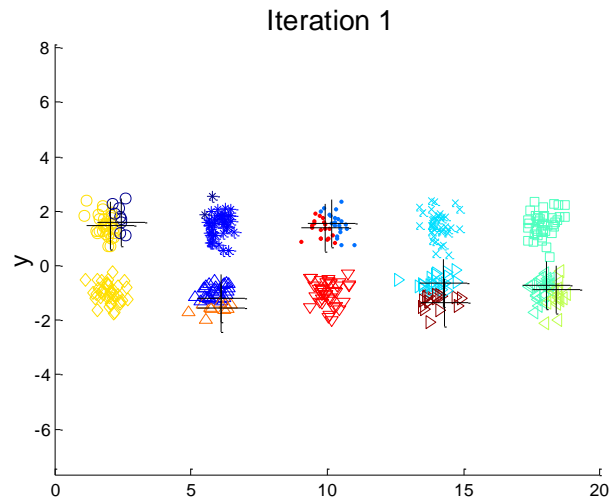
- Depending on the choice of initial centroids, B and C may get merged or remain separate

# 10 Clusters Example



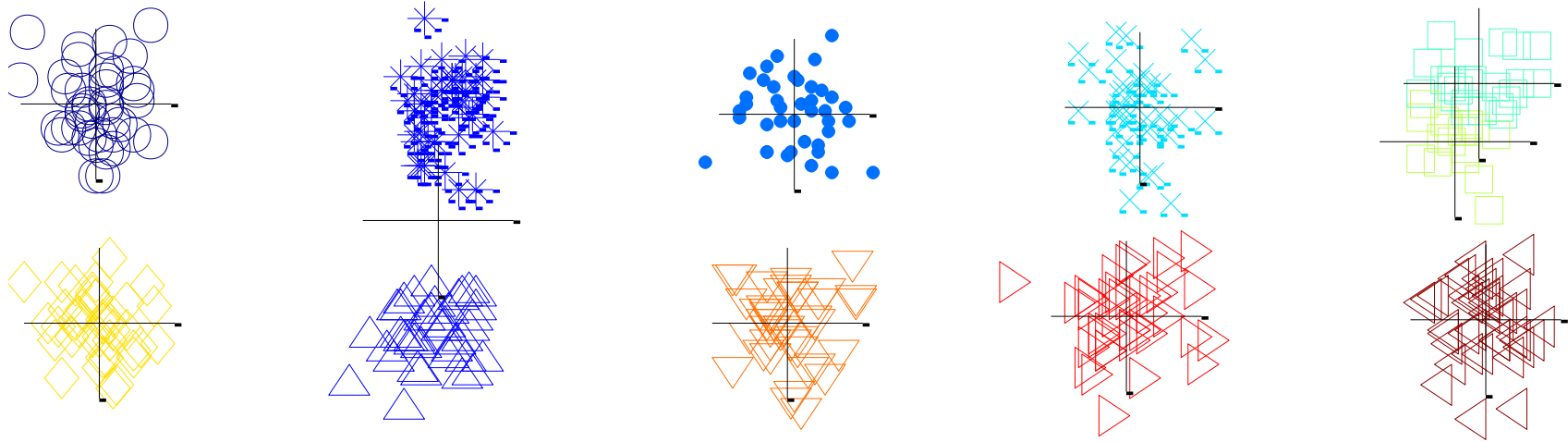
**Starting with two initial centroids in one cluster of each pair of clusters**

# 10 Clusters Example



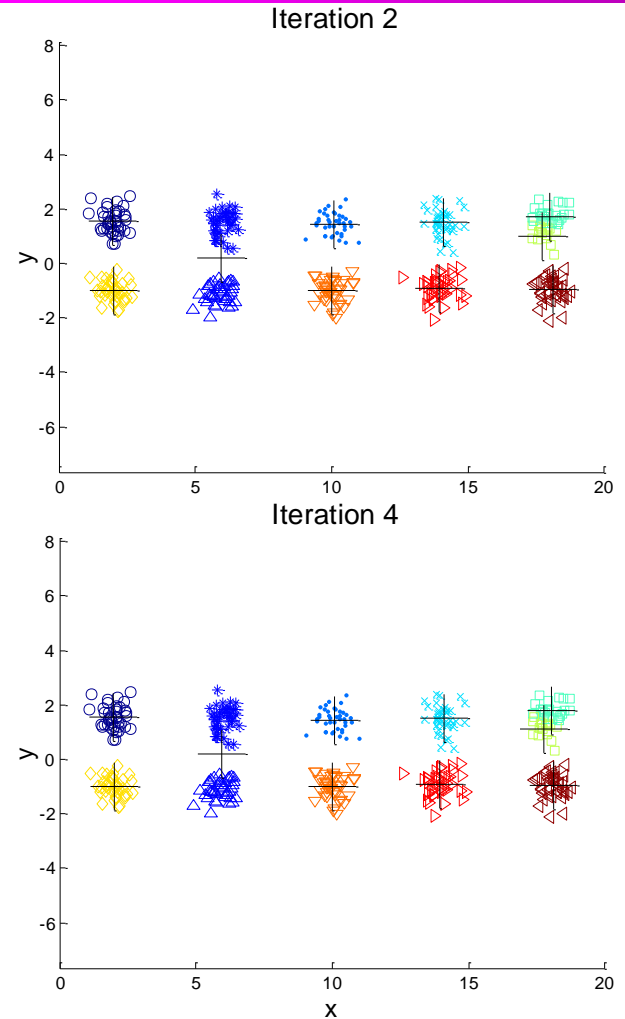
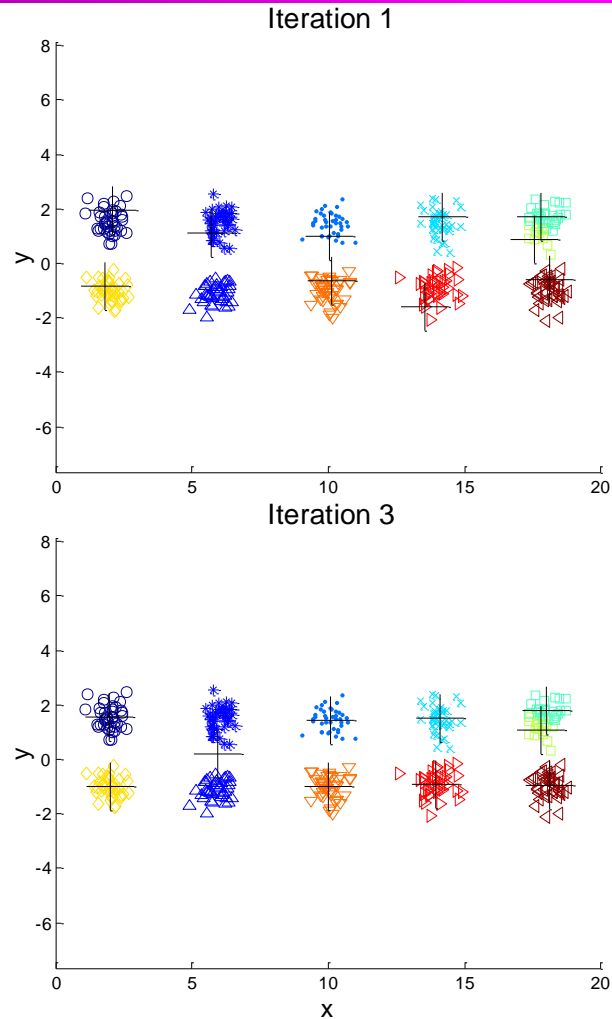
**Starting with two initial centroids in one cluster of each pair of clusters**

# 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
  - Use hierarchical clustering to determine initial centroids

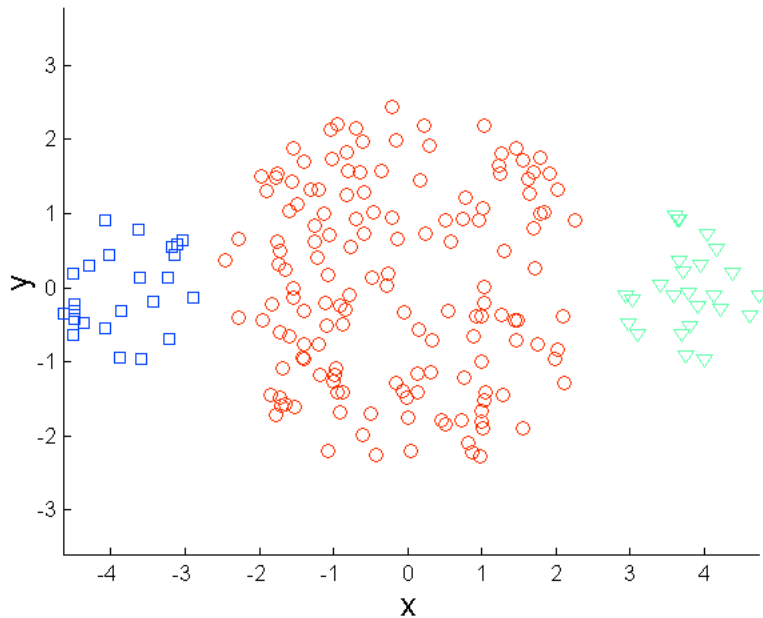
# Limitations of K-means

---

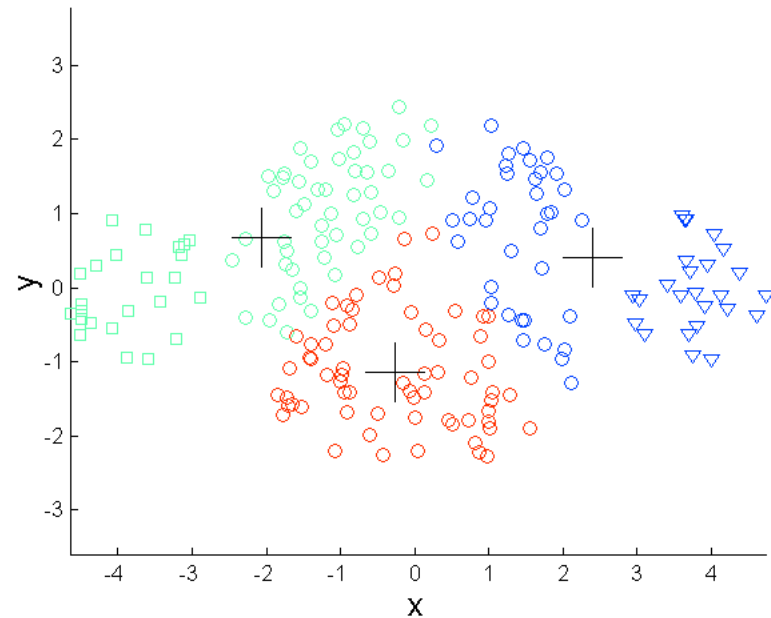
- K-means has problems when clusters are of differing
  - Sizes
  - Densities
  - Non-globular shapes
- K-means has problems when the data contains outliers.
  - One possible solution is to remove outliers before clustering



# Limitations of K-means: Differing Sizes

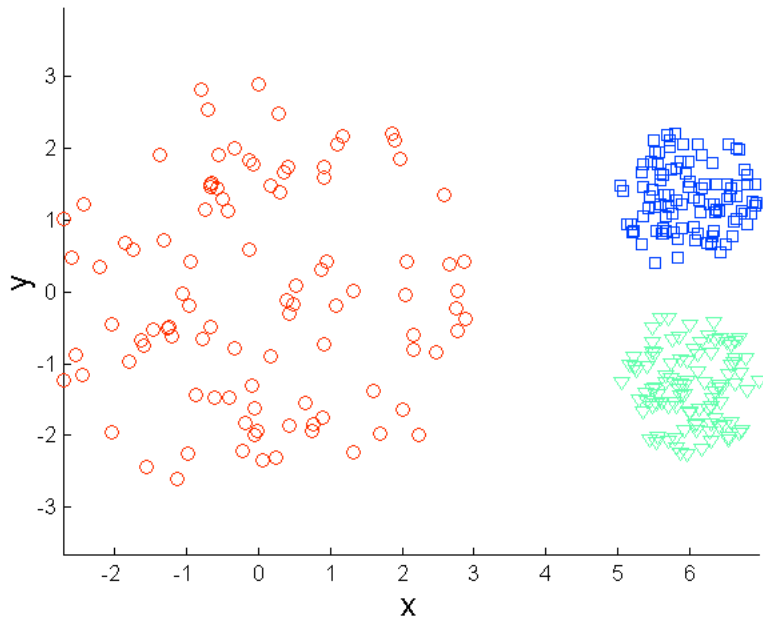


**Original Points**

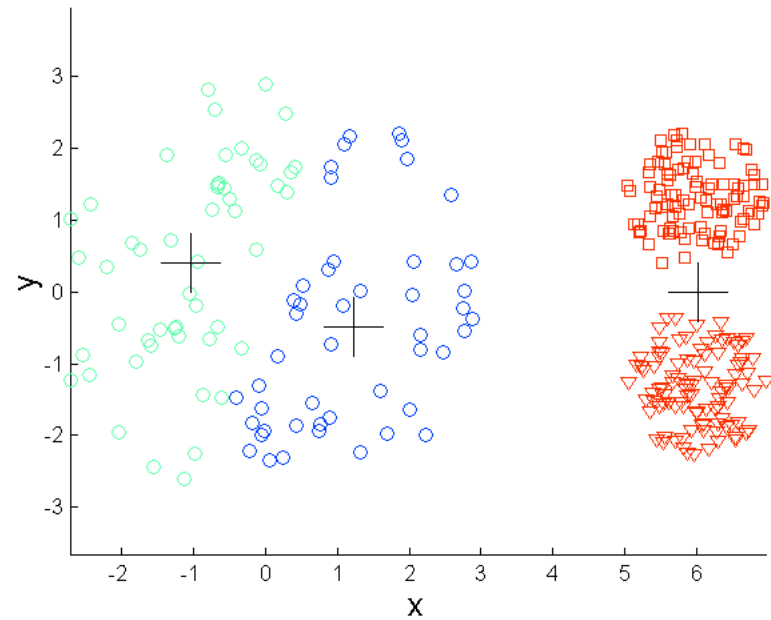


**K-means (3 Clusters)**

# Limitations of K-means: Differing Density

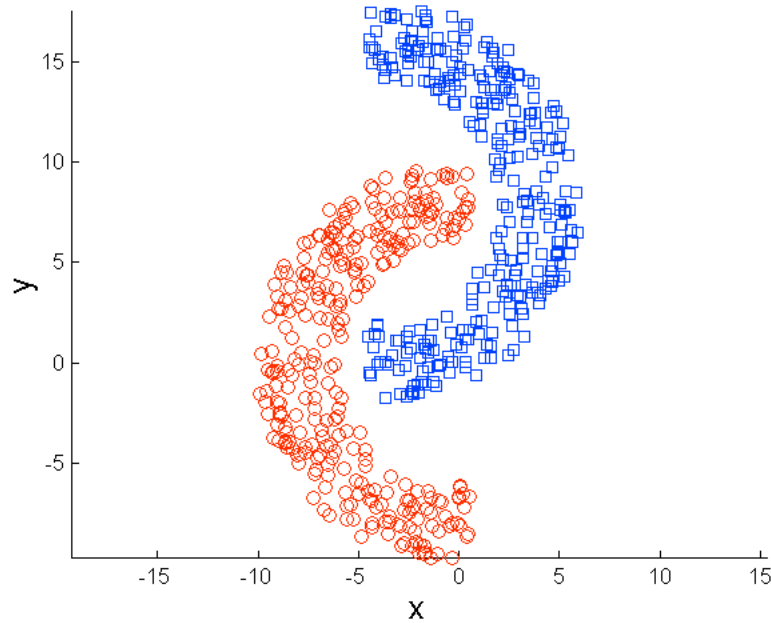


**Original Points**

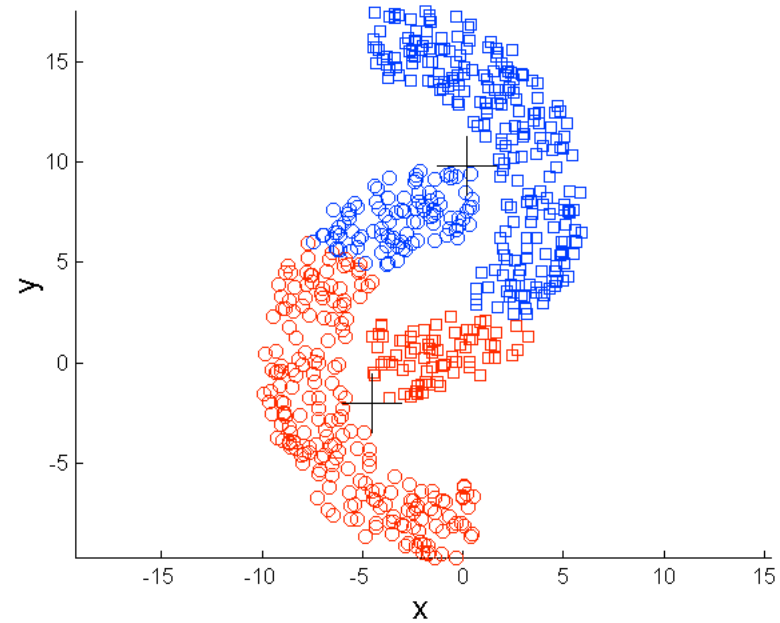


**K-means (3 Clusters)**

# Limitations of K-means: Non-globular Shapes

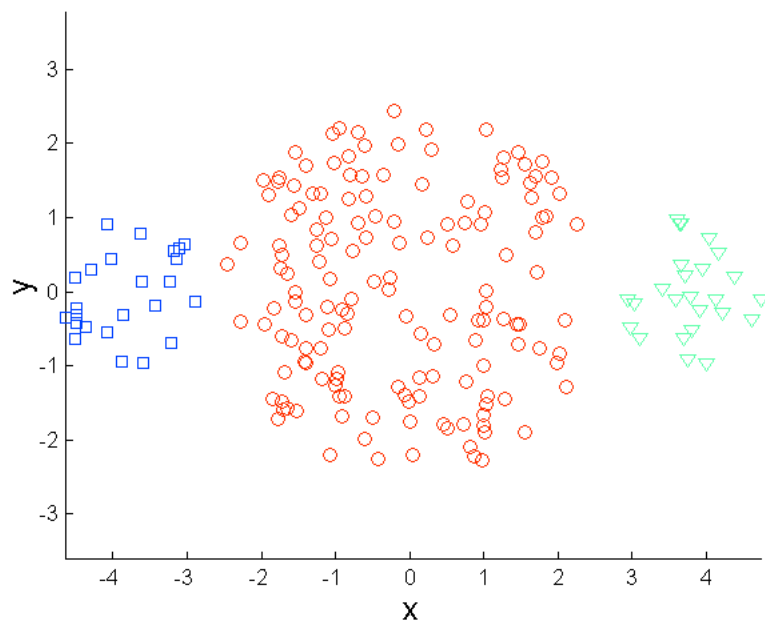


**Original Points**

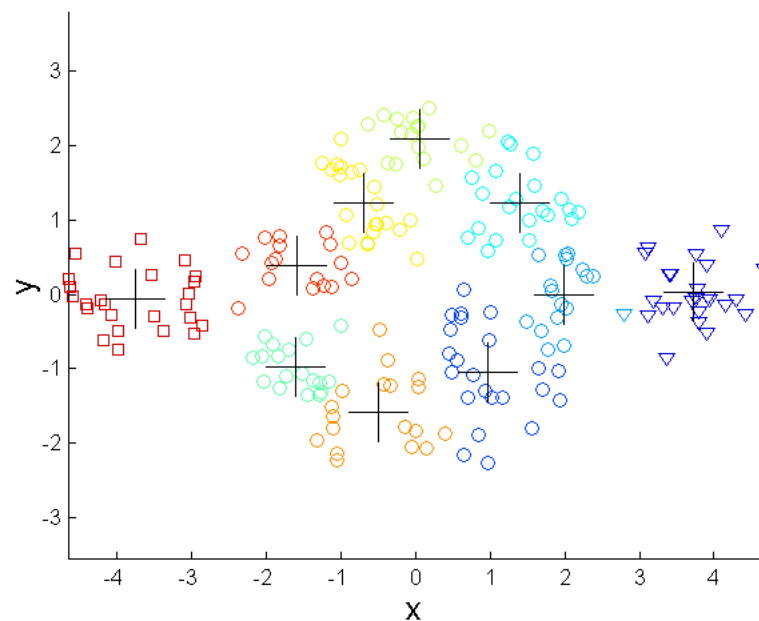


**K-means (2 Clusters)**

# Overcoming K-means Limitations



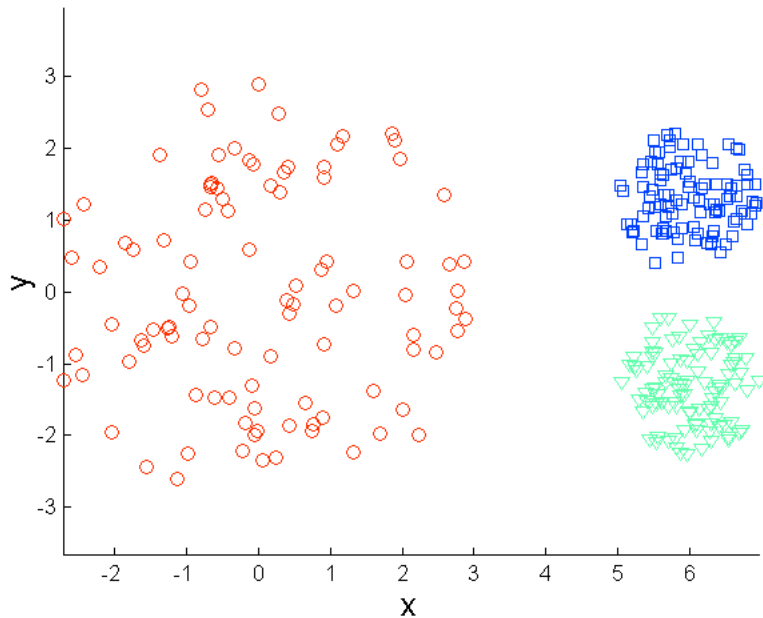
**Original Points**



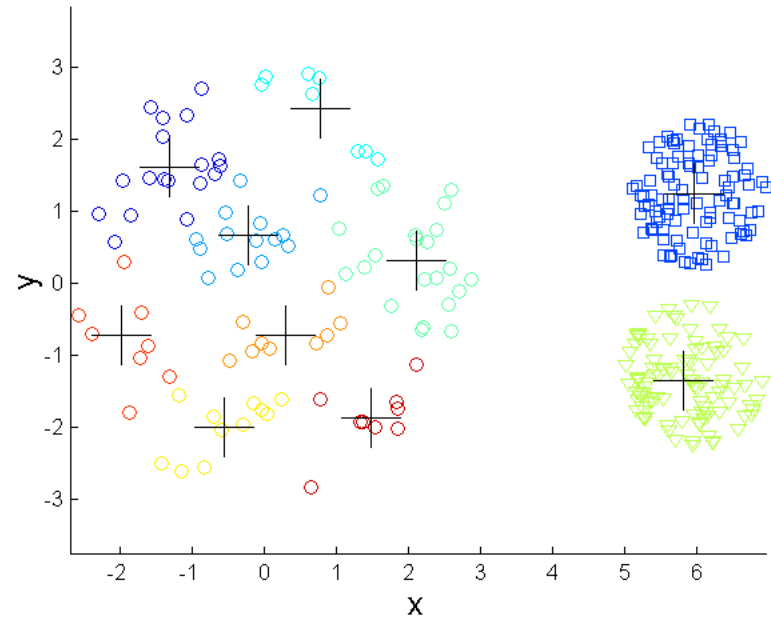
**K-means Clusters**

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



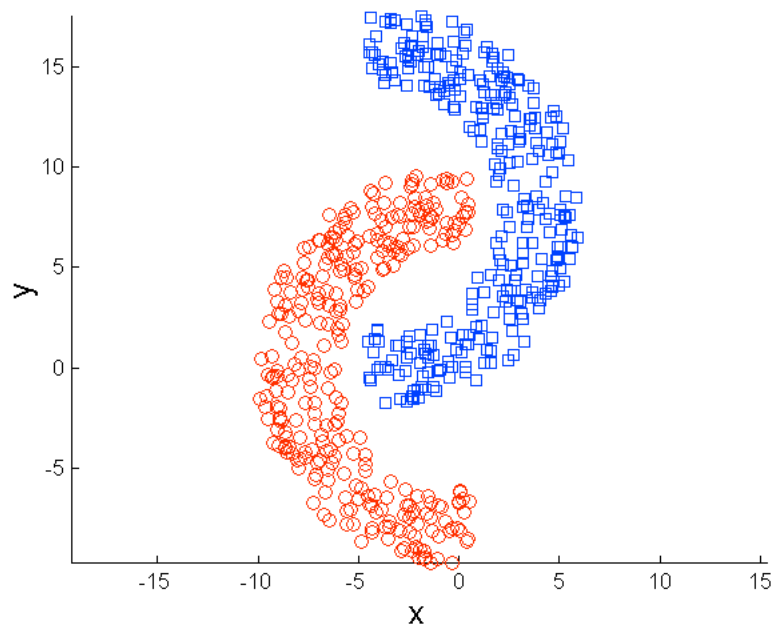
**Original Points**



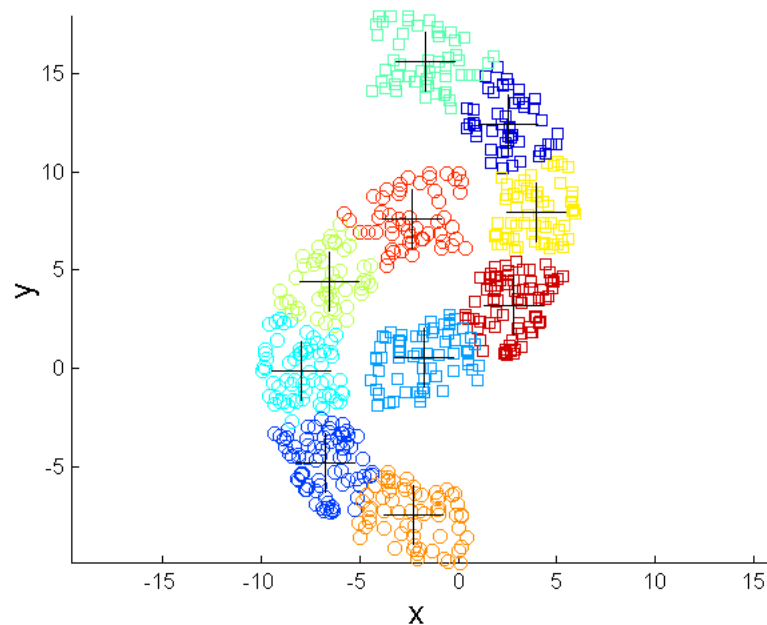
**K-means Clusters**

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.



# Questions

# Density Based Clustering

---

- Clusters are regions of high density that are separated from one another by regions of low density.





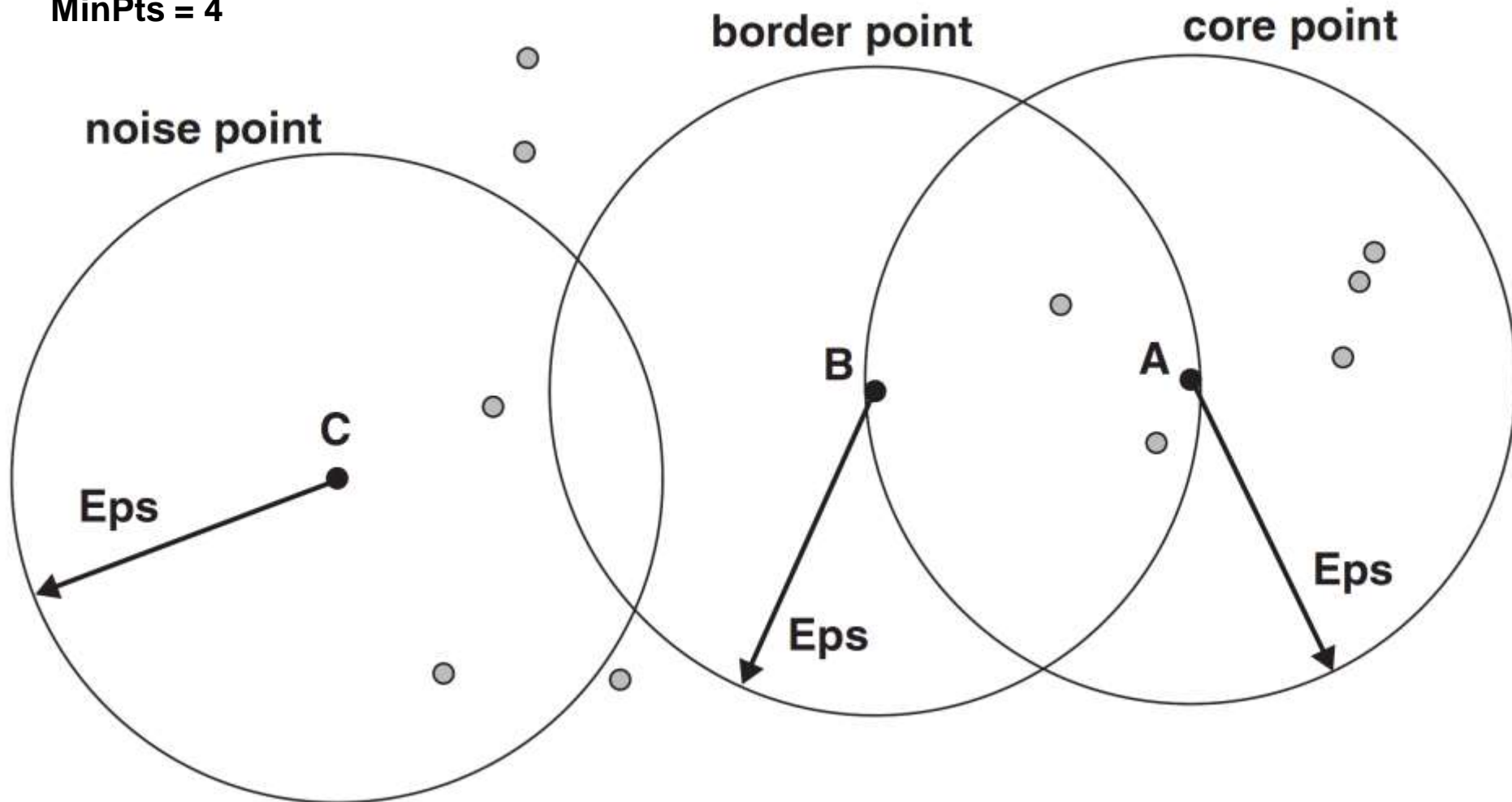
# DBSCAN

---

- 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, and Noise Points

MinPts = 4

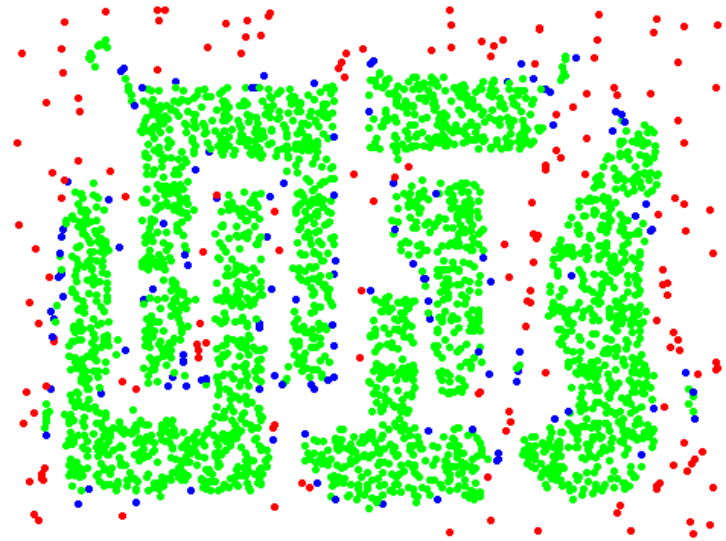


# DBSCAN: Core, Border and Noise Points

---



Original Points



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

**Eps = 10, MinPts = 4**

# DBSCAN Algorithm

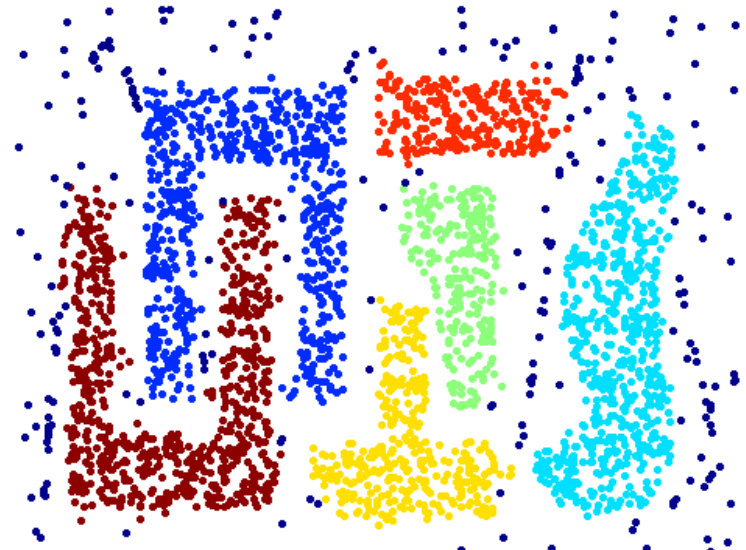
---

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

# When DBSCAN Works Well



Original Points



Clusters (dark blue points indicate noise)

- Can handle clusters of different shapes and sizes
- Resistant to noise



# Questions