# CS722/822: Machine Learning
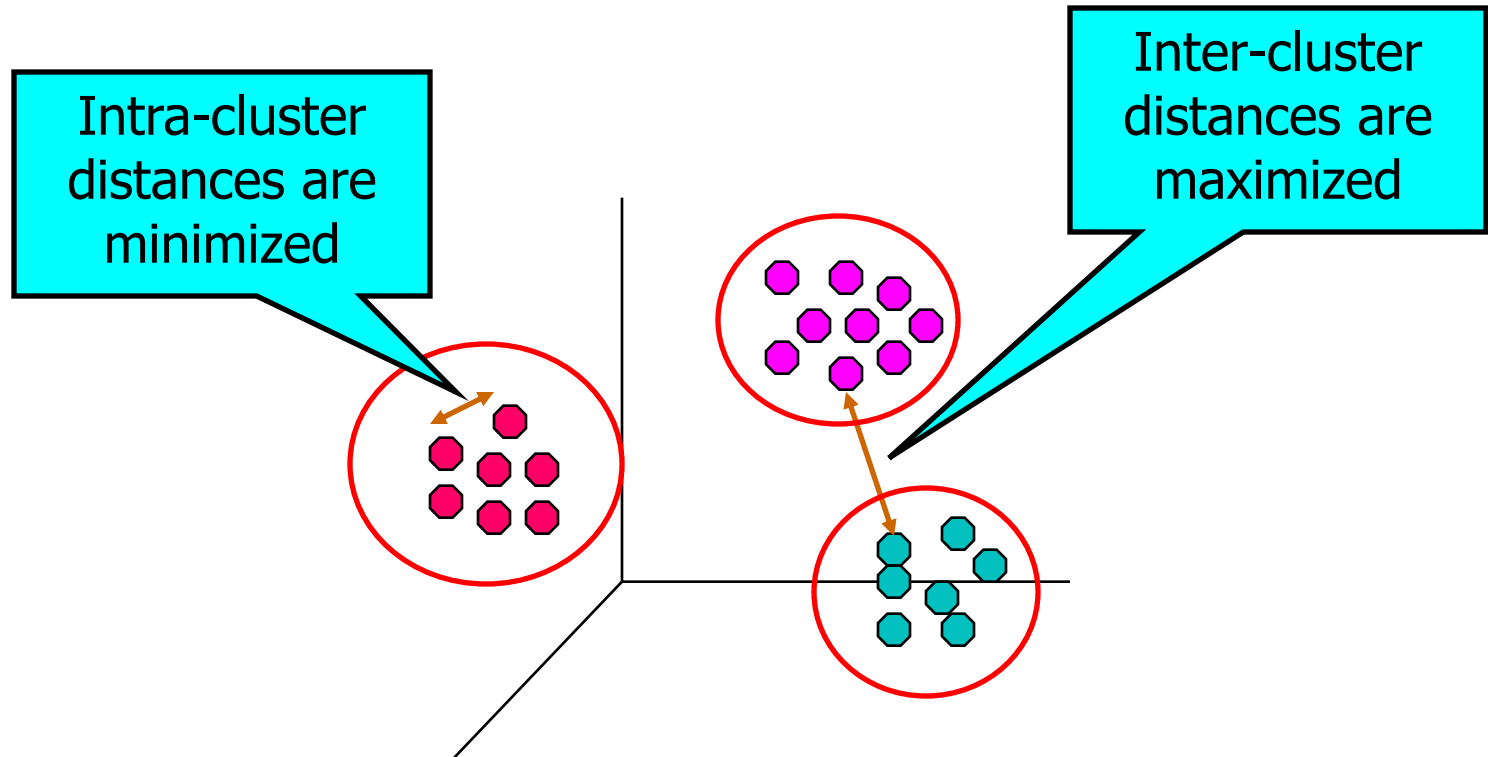
Instructor: Jiangwen Sun

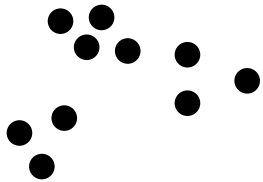Computer Science Department

# Unsupervised Learning

● Draw inferences from the data for exploratory analysis, such as finding hidden patterns or grouping of examples

- – Given a collection of examples, each consisting of a set of features, but without a clear target, i.e., $\{x_1, x_2, \cdots, x_N\}$

- – Examples:

  - ◆ Principal component analysis (PCA)

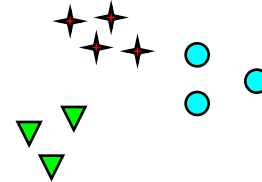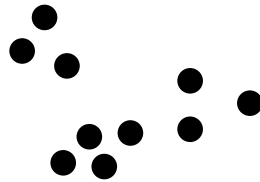  - ◆ Clustering (cluster analysis)

# 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

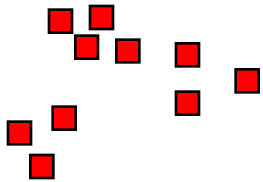Intra-cluster distances are minimized

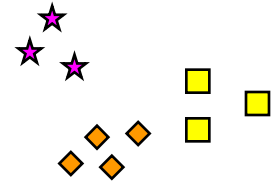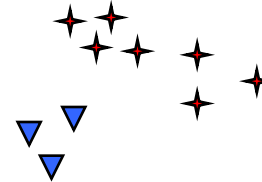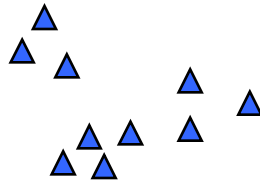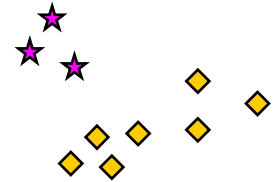Inter-cluster distances are maximized

# Notion of a Cluster can be Ambiguous

How many clusters?

Six Clusters

Two Clusters
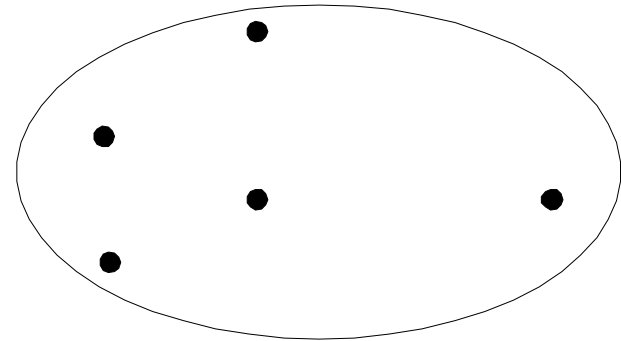
Four Clusters

# Types of Clusterings

- A clustering is a set of clusters

- Important distinction between hierarchical and partitional sets of clusters

- Partitional Clustering
  - A division data objects into **non-overlapping** subsets (clusters) such that each data object is in exactly **one** subset

- Hierarchical clustering
  - A set of nested clusters organized as a hierarchical tree

# Partitional Clustering



**Original Points**                              **A Partitional  Clustering**

# Hierarchical Clustering



Nested Clusters

Dendrogram

# Other Distinctions Between Sets of Clusters

- Non-exclusive versus exclusive
  - In non-exclusive clustering, points may belong to multiple clusters.
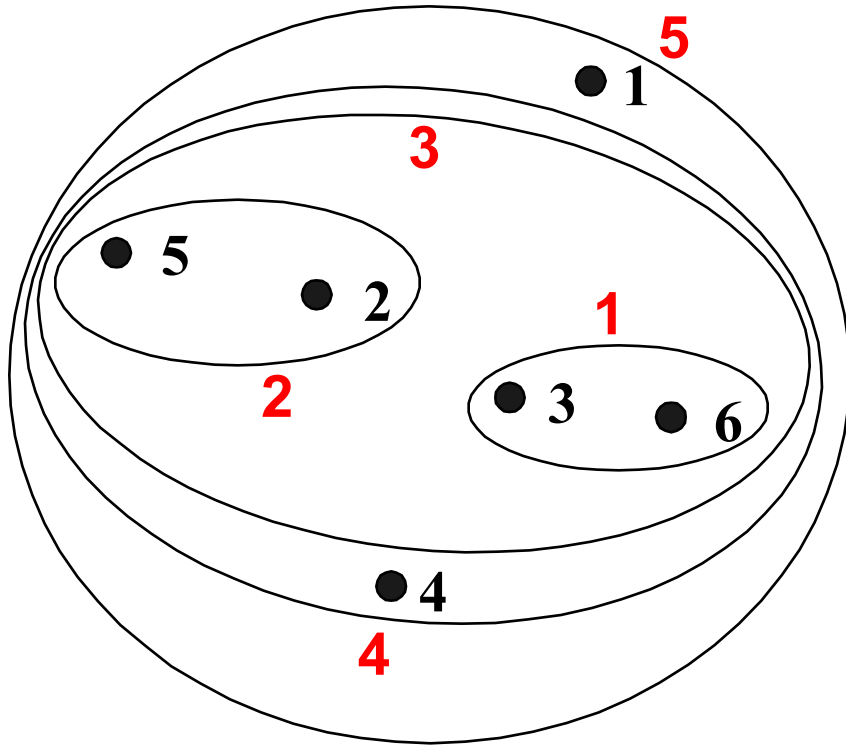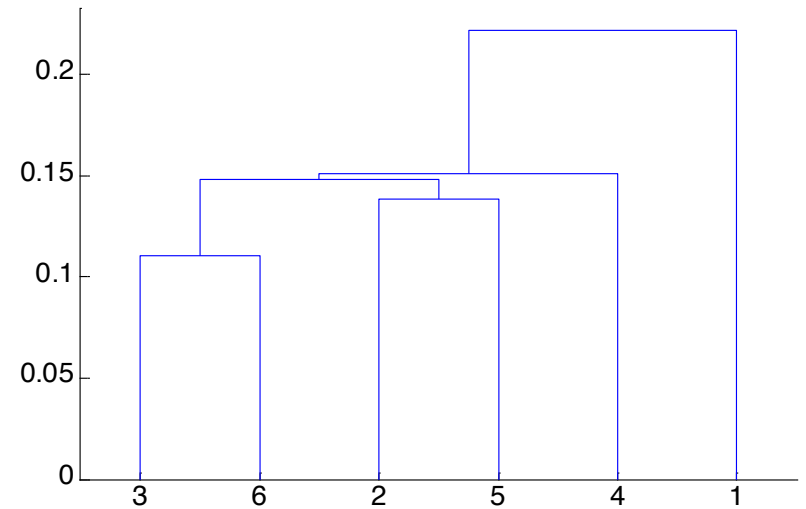  - Can represent multiple classes or 'border' points
- Fuzzy versus non-fuzzy, probability vs non-probability
  - In fuzzy clustering, a point belongs to every cluster with some weights
  - Often with constraints that weights must sum to 1
  - Probabilistic clustering has similar characteristics
- Partial 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

- Well-separated clusters

- Center-based clusters

- Contiguous clusters

- Density-based clusters

- Property or Conceptual

- 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.
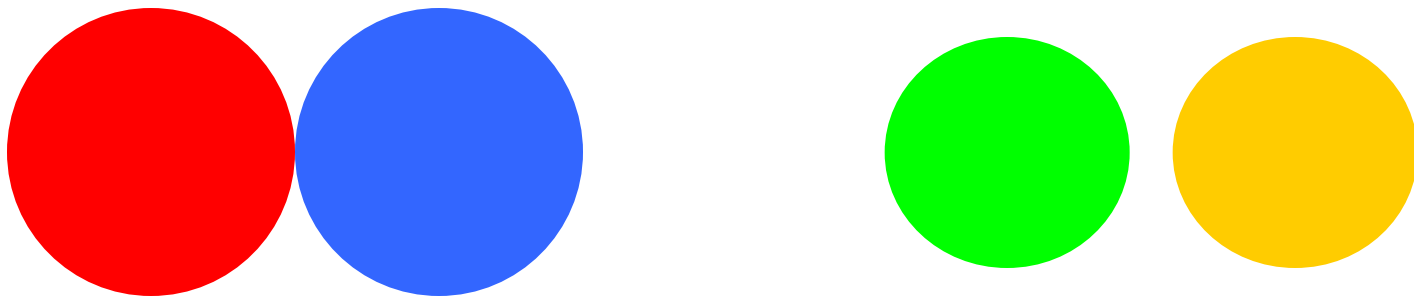
3 well-separated clusters
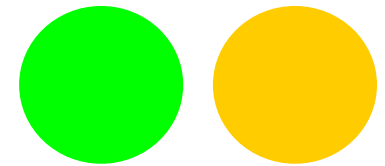
# Types of Clusters: Center-Based

● Center-based

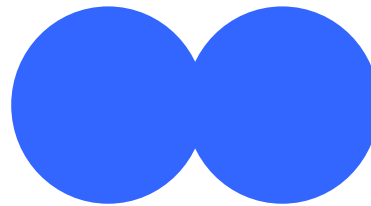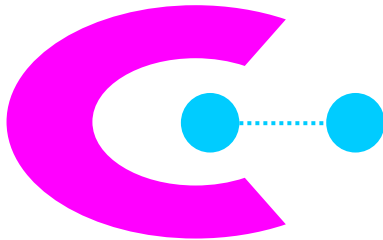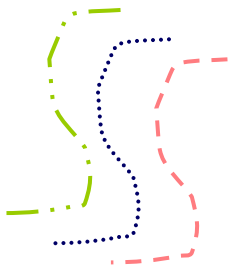– A cluster is a set of points such that any point in a cluster is closer (more similar) to the "center" of a cluster, than to the center of any other cluster

– The center of a cluster is often the <span style="color:red">centroid</span>, the average of all the points in the cluster, or the <span style="color:red">medoid</span>, the most "representative" point of a cluster

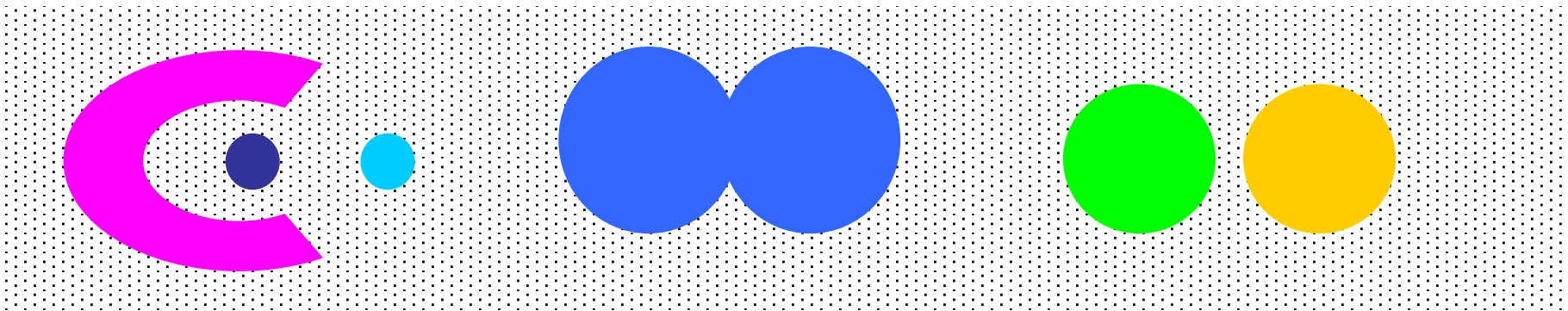4 center-based clusters

# 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.

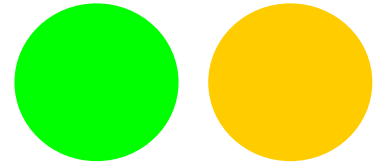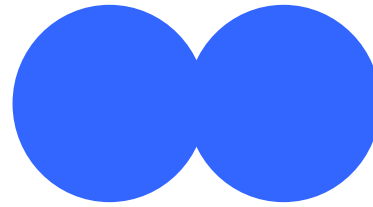8 contiguous clusters

# 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
    - noise and outliers are present

6 density-based clusters

# Contiguity Based vs Density-Based

8 contiguous clusters

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.

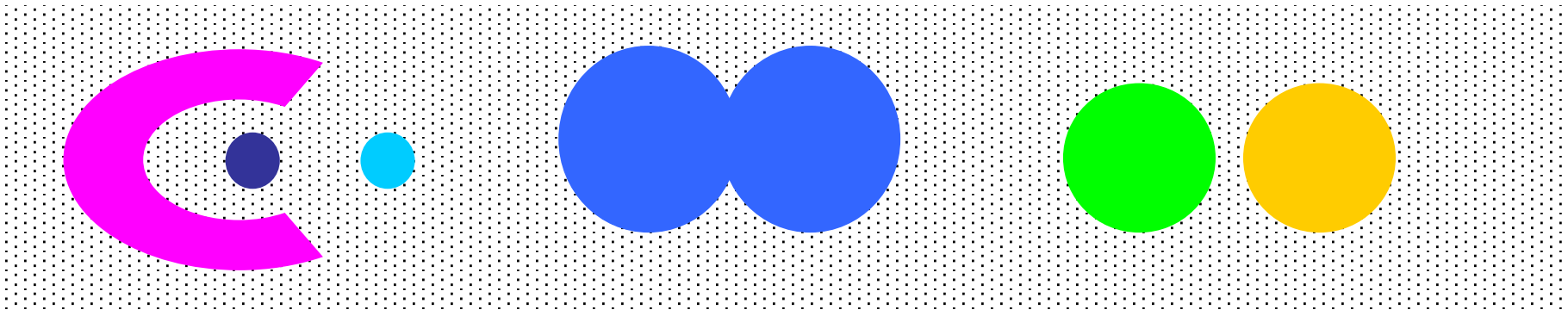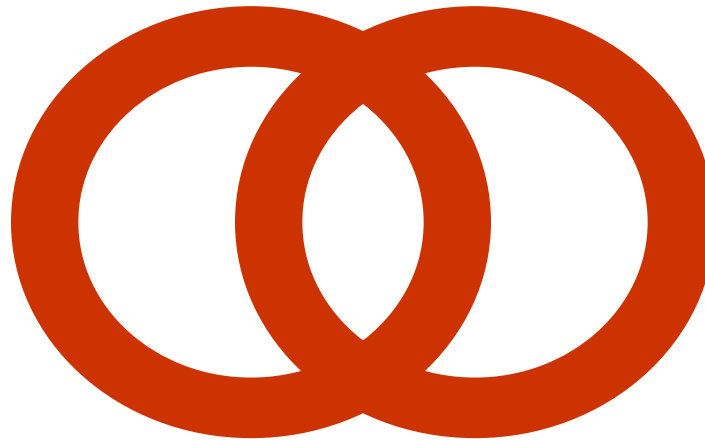- Enumerate all possible ways of dividing the points into clusters and evaluate the `goodness' of each potential set of clusters by using the given objective function.  (NP Hard)

- 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 (variable) type
  - Dictates type of similarity
- Type of Data
  - Dictates type of similarity
- Dimensionality
- Noise and Outliers
- Type of Distribution