



COMP7/8118 M50

# Data Mining

Clustering Concepts

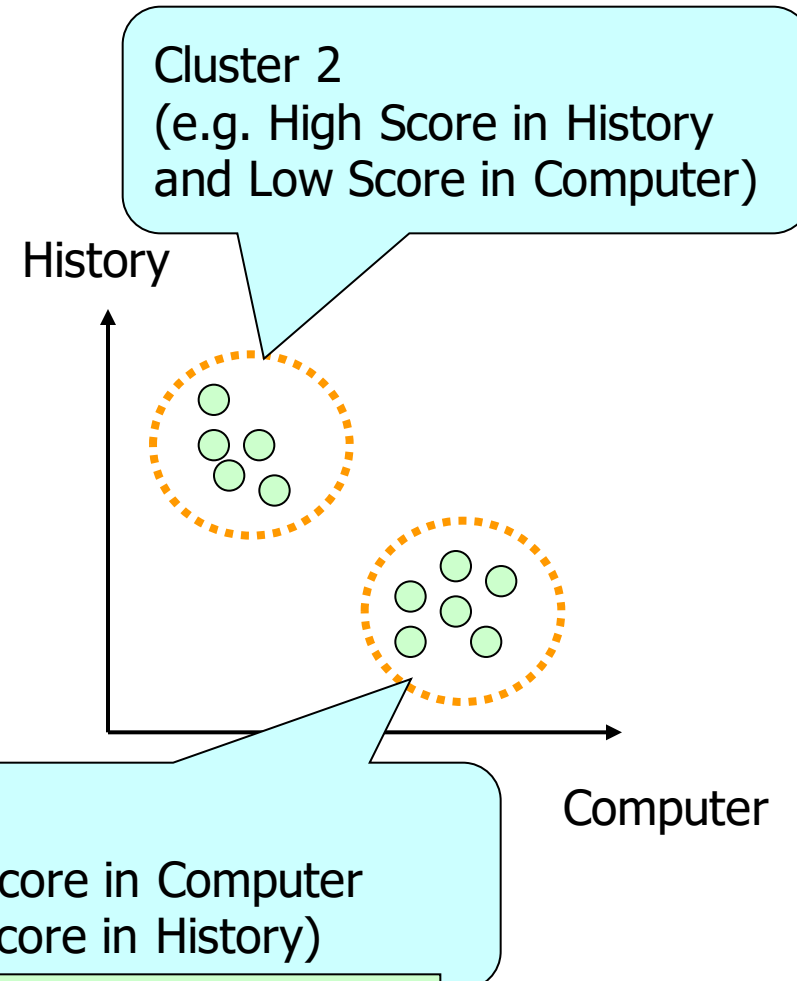
Xiaofei Zhang

*Slides compiled from Jiawei Han and Raymond C.W. Wong's work*

THE UNIVERSITY OF  
**MEMPHIS**

# Clustering

	Computer	History
Raymond	100	40
Louis	90	45
Wyman	20	95
...	...	...



Problem: to find all clusters

# Why Clustering?

- **Understanding**

- **Group** related **documents** for browsing, **genes and proteins** that have similar functionality, **stocks** with similar price fluctuations, users with same behavior

- **Summarization**

- Reduce the size of large data sets

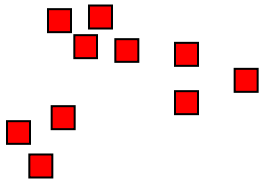
- **Applications**

- Biology - Group different species
- Psychology and Medicine - Group medicine
- Business - Group different customers for marketing
- Network - Group different types of traffic patterns
- Software - Group different programs for data analysis

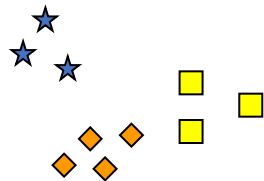
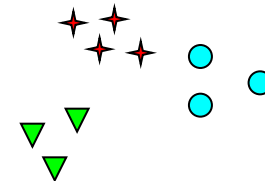
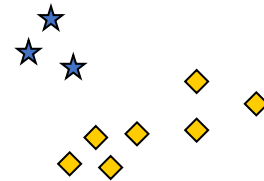
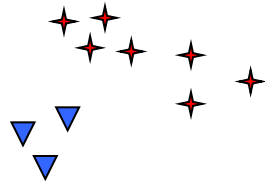
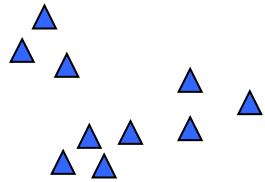
# Notion of a Cluster can be Ambiguous



How many clusters?



Two Clusters



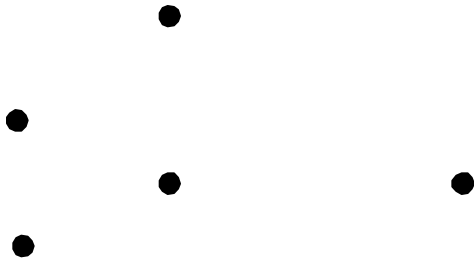
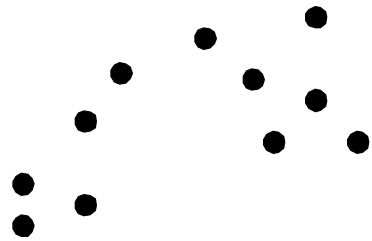
Four Clusters

Six Clusters

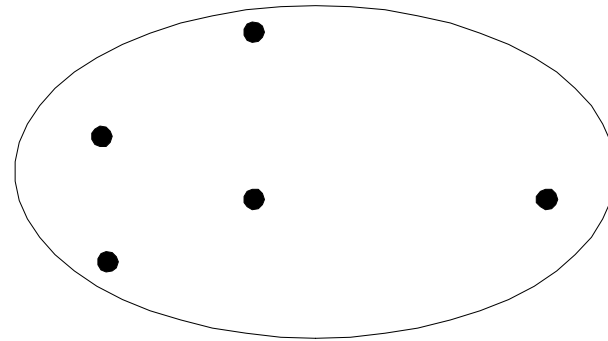
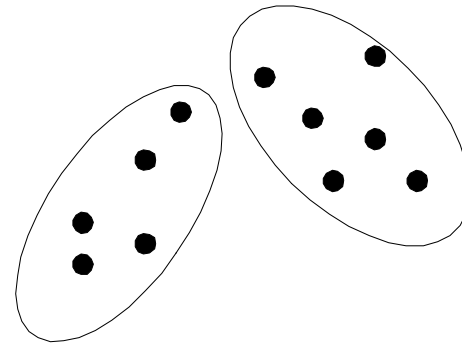
# Types of Clustering

- A **clustering** is a set of **clusters**
- Important distinction between **hierarchical** and **partitional** sets of clusters
- **Partitional** Clustering
  - A division of data objects into 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

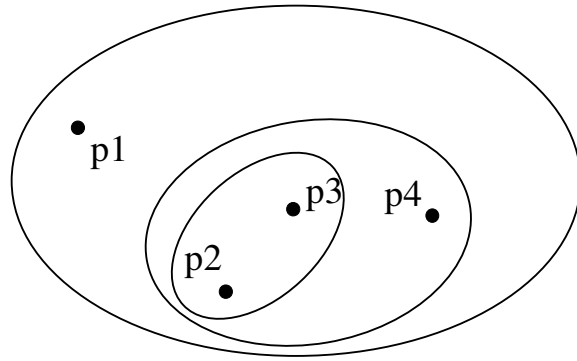


Original Points

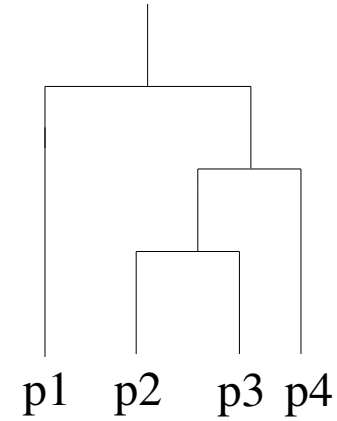


A Partitional Clustering

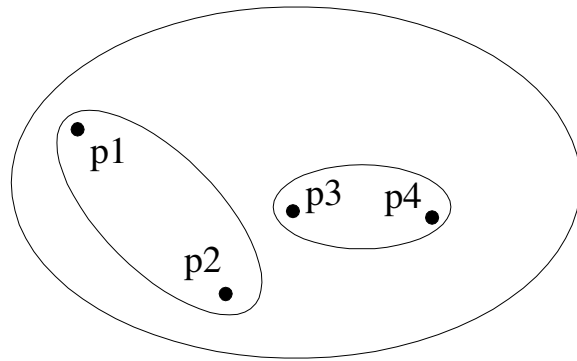
# Hierarchical Clustering



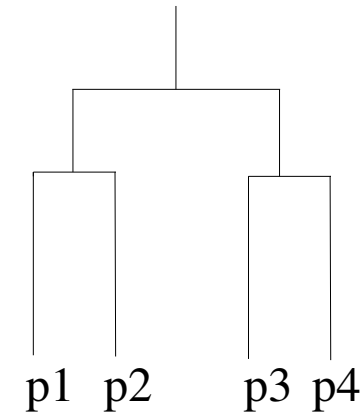
Traditional Hierarchical Clustering



Traditional Dendrogram



Non-traditional Hierarchical Clustering



Non-traditional Dendrogram

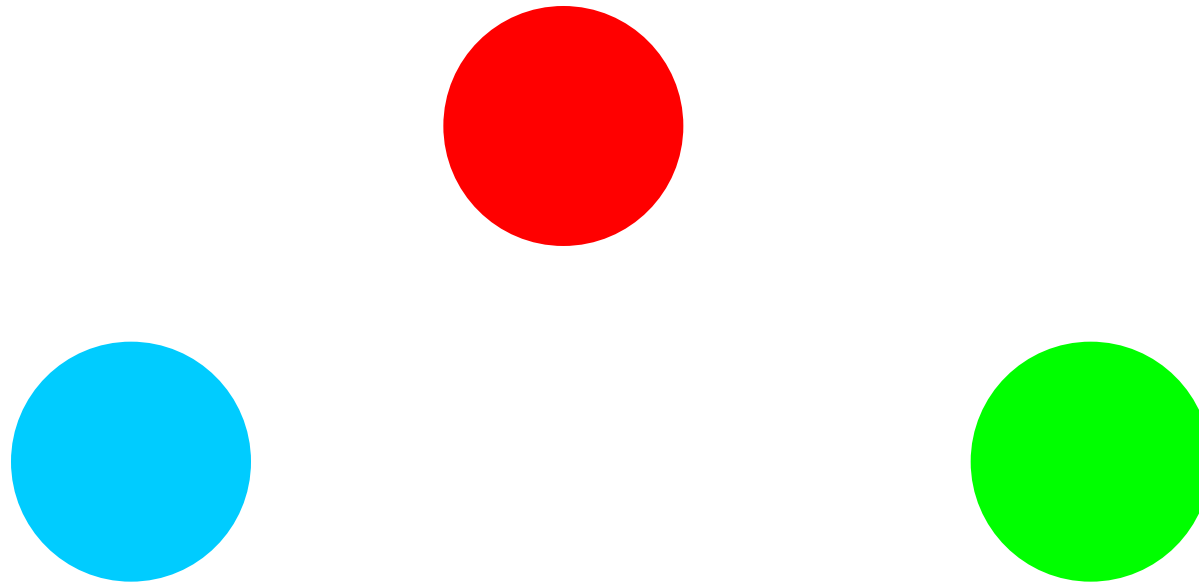
# Other types of clustering

- **Exclusive** (or **non-overlapping**) versus **non-exclusive** (or **overlapping**)
  - In non-exclusive clustering, points may belong to multiple clusters.
    - Points that belong to multiple classes, or 'border' points
- **Fuzzy** (or **soft**) versus **non-fuzzy** (or **hard**)
  - In fuzzy clustering, a point belongs to every cluster with some weight between 0 and 1
    - Weights usually must sum to 1 (often interpreted as **probabilities**)



# Types of Clusters: Well-Separated

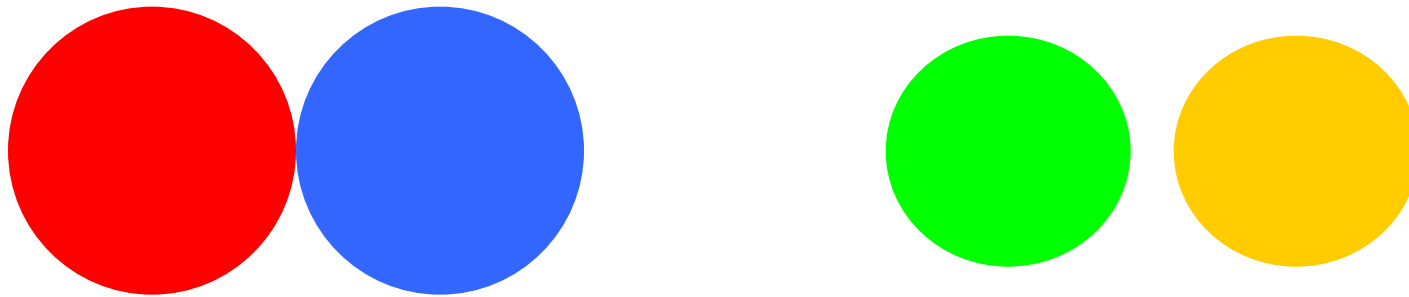
- **Well-Separated Clusters:**
  - A cluster is a set of points such that any point in a cluster is closer (or more similar) to every other point in the cluster than to any point not in the cluster.



3 well-separated clusters

# Types of Clusters: Center-Based

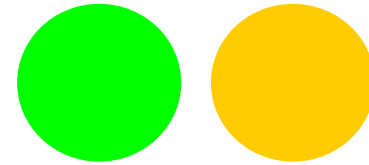
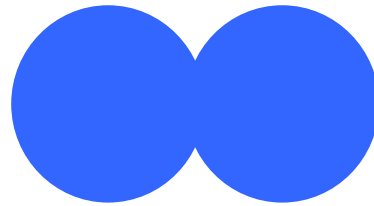
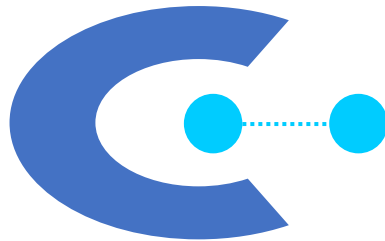
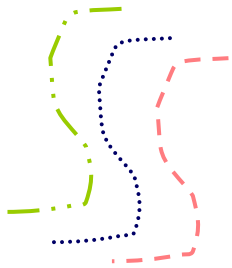
- Center-based
  - A cluster is a set of objects such that an object in a cluster is **closer** (more **similar**) to the “center” of a cluster, than to the center of any other cluster
  - The center of a cluster is often a **centroid**, the minimizer of distances from all the points in the cluster,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



4 center-based clusters

# Types of Clusters: Contiguity-Based

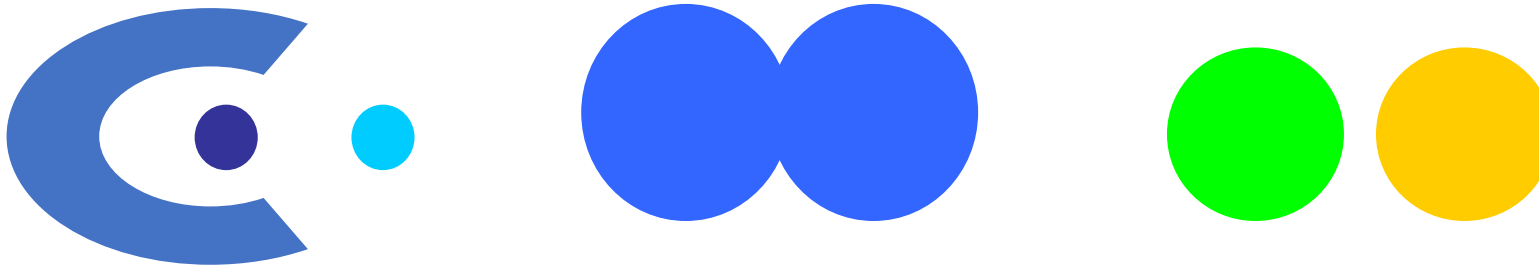
- Contiguous Cluster (Nearest neighbor or Transitive)
  - A cluster is a set of points such that a point in a cluster is closer (or more similar) to one or more other points in the cluster than to any point not in the cluster.



8 contiguous clusters

# Types of Clusters: Density-Based

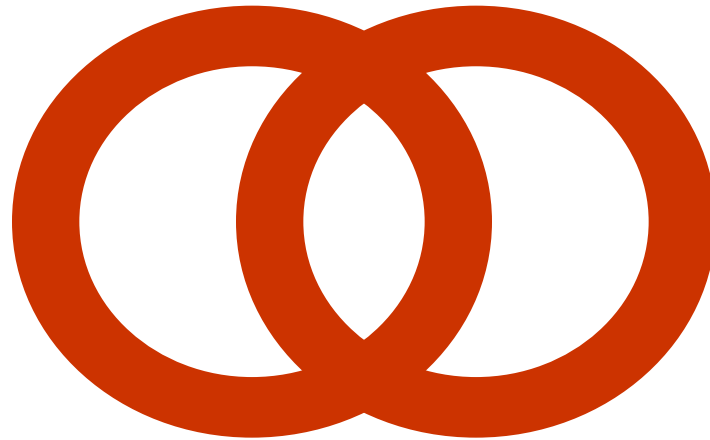
- Density-based
  - A cluster is a dense region of points, which is separated by low-density regions, from other regions of high density.
  - Used when the clusters are irregular or intertwined, and when noise and outliers are present.



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

# Objective Function

- Clustering as an **optimization problem**
  - 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**.
    - The **parameters** for the model are determined from the data, and they determine the clustering
    - E.g., **Mixture models** assume that the data is a 'mixture' of a number of statistical distributions.

# Takeaways

- The notion of clustering can be ambiguous
- Clustering results can be very different under different semantics or using different methods
- The evaluation of a clustering output is non-trivial (will be covered later)