

# Data Mining -- Clustering

Instructor: Jen-Wei Huang

Office: 92528 in the EE building  
jwhuang@mail.ncku

## Clustering

- ▶ Cluster: a collection of data objects
  - Similar to one another within the same cluster
  - Dissimilar to the objects in other clusters
- ▶ Cluster analysis
  - Finding similarities between data according to the characteristics found in the data and grouping similar data objects into clusters
- ▶ **Unsupervised learning**: no predefined classes

# Typical Applications

- ▶ As a **stand-alone tool** to get insight into data distribution or as a **preprocessing step** for other algorithms
- ▶ Pattern Recognition
- ▶ Image Processing
- ▶ Economic Science (especially market research)
- ▶ WWW
  - Web pages (resources) clustering
  - Cluster Weblog data to discover groups of similar access patterns

## Examples

- ▶ Marketing: Help marketers discover distinct groups in their customer bases, and then use this knowledge to develop targeted marketing programs
- ▶ Land use: Identification of areas of similar land use in an earth observation database
- ▶ Insurance: Identifying groups of motor insurance policy holders with a high average claim cost
- ▶ City-planning: Identifying groups of houses according to their house type, value, and geographical location

# Quality of Clustering

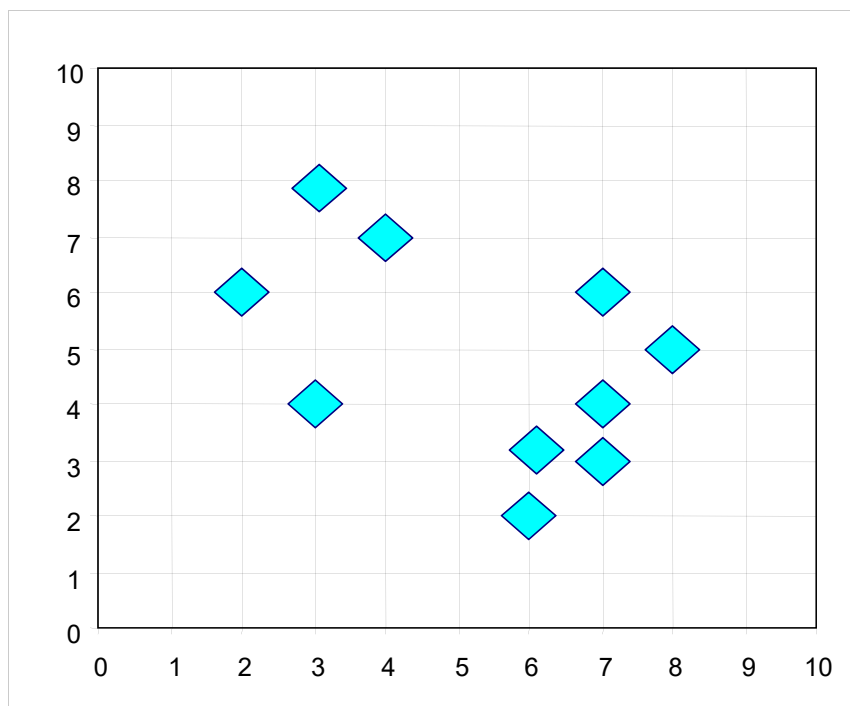
- ▶ Good clusters :
  - high intra-class similarity
  - low inter-class similarity
- ▶ The quality of a clustering result depends on both the similarity measure used by the method and its implementation
- ▶ The quality of a clustering method is also measured by its ability to discover some or all of the hidden patterns

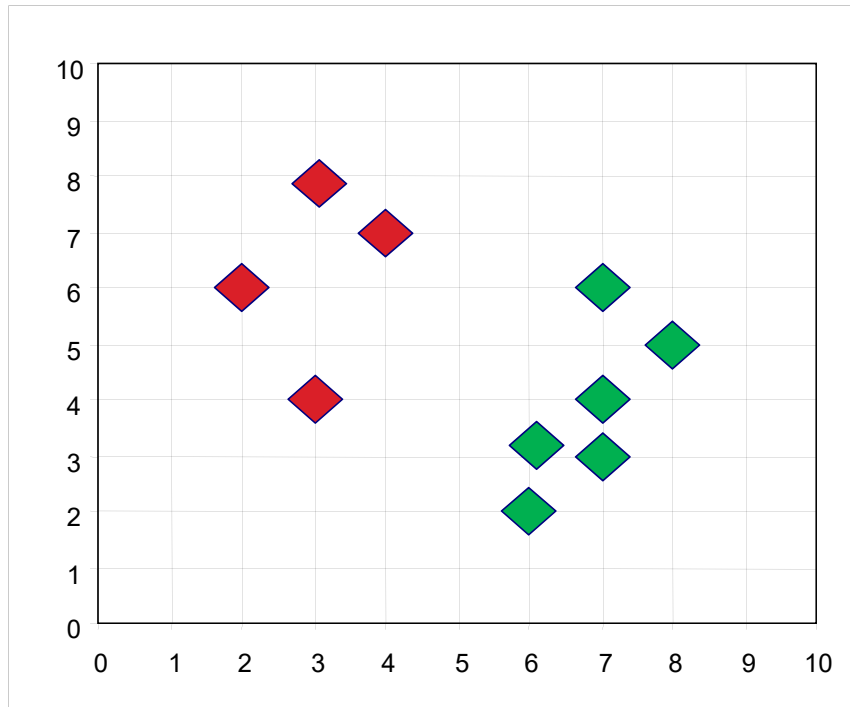
## Measure the Quality

- ▶ **Dissimilarity/Similarity metric**: Similarity is expressed in terms of a distance function, typically metric:  $d(i, j)$
- ▶ There is a separate “quality” function that measures the “goodness” of a cluster.
- ▶ The definitions of **distance functions** are usually very different for interval-scaled, boolean, categorical, ordinal ratio, and vector variables.
- ▶ Weights should be associated with different variables based on applications and data semantics.
- ▶ It is hard to define “similar enough” or “good enough”
  - the answer is typically highly subjective.

# Requirements of Clustering

- ▶ Able to deal with noise and outliers
- ▶ Insensitive to order of input records
- ▶ High dimensionality
- ▶ Incorporation of user-specified constraints
- ▶ Interpretability and usability
- ▶ Scalability
- ▶ Ability to deal with different types of attributes
- ▶ Ability to handle dynamic data
- ▶ Discovery of clusters with arbitrary shape
- ▶ Determination input parameters





## References

- ▶ Slides from Prof. J.-W. Han, UIUC
- ▶ Slides from Prof. M.-S. Chen, NTU
- ▶ Slides from Prof. W.-Z. Peng, NCTU