

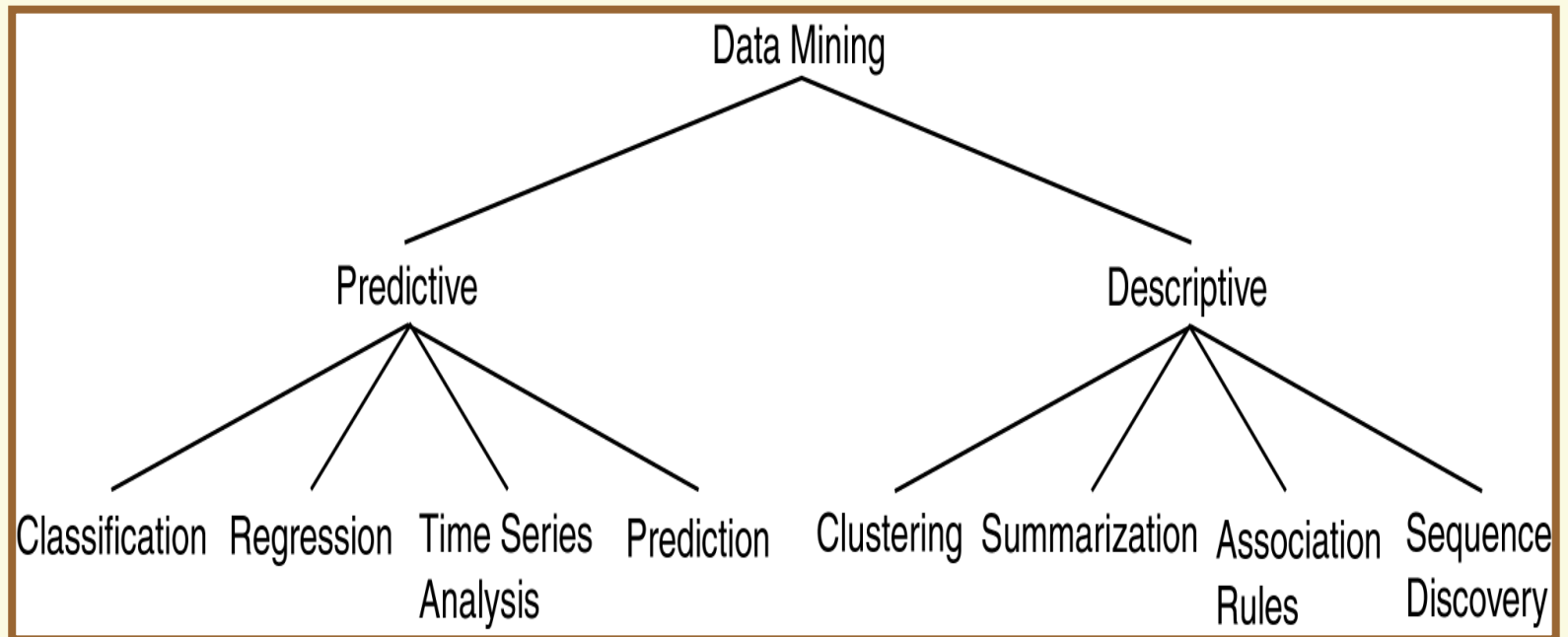
CPT_S 415

Big Data

Clustering

- ✓ Data clustering
- ✓ Graph clustering

Data Mining Models and Tasks



Use variables to predict unknown or future values of other variables.

Find human-interpretable patterns that describe the data.

Data clustering


What is Cluster Analysis?

- ✓ Cluster: A collection of data objects
 - similar (or related) to one another within the same group
 - dissimilar (or unrelated) to the objects in other groups
- ✓ Cluster analysis (or *clustering*, *data segmentation*, ...)
 - Finding similarities between data according to the characteristics found in the data and grouping similar data objects into clusters
- ✓ **Unsupervised learning**: no predefined classes (i.e., *learning by observations* vs. learning by examples: supervised)
- ✓ Typical applications
 - As a **stand-alone tool** to get insight into data distribution
 - As a **preprocessing step** for other algorithms

Applications of Cluster Analysis

- ✓ Data reduction
 - Summarization: Preprocessing for regression, classification, and association analysis
 - Compression: Image processing: vector quantization
- ✓ Prediction based on groups
 - Cluster & find characteristics/patterns for each group
- ✓ Finding K-nearest Neighbors
 - Localizing search to one or a small number of clusters
- ✓ Outlier detection: Outliers are often viewed as those “far away” from any cluster

Clustering: Application Examples

- ✓ Biology: taxonomy of living things: kingdom, phylum, class, order, family, genus and species
- ✓ Land use: Identification of areas of similar land use in an earth observation database
- ✓ Marketing: Help marketers discover distinct groups in their customer bases, and then use this knowledge to develop targeted marketing programs
- ✓ City-planning: Identifying groups of houses according to their house type, value, and geographical location
- ✓ Earth-quake studies: Observed earth quake epicenters should be clustered along continent faults
- ✓ Climate: understanding earth climate, find patterns of atmospheric and ocean
- ✓ Traffic flow analyses 

Basic Steps to Develop a Clustering Task

- ✓ Feature selection
 - Select info concerning the task of interest
 - Minimal information redundancy
- ✓ Proximity measure
 - Similarity of two feature vectors
- ✓ Clustering criterion
 - Expressed via a cost function or some rules
- ✓ Clustering algorithms
 - Choice of algorithms
- ✓ Validation of the results
 - Validation test (also, *clustering tendency* test)
- ✓ Interpretation of the results
 - Integration with applications

$$H_{ind} = \frac{\sum_{i=1}^n q_i}{\sum_{i=1}^n q_i + \sum_{i=1}^n w_i}$$

What Is Good Clustering?

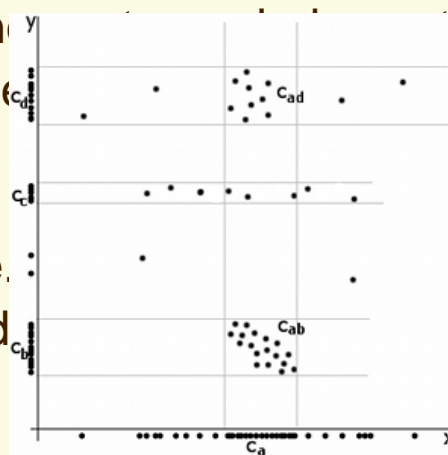
- ✓ A good clustering method will produce high quality clusters
 - high intra-class similarity: **cohesive** within clusters
 - low inter-class similarity: **distinctive** between clusters
- ✓ The quality of a clustering method depends on
 - the similarity measure used by the method
 - its implementation, and
 - Its ability to discover some or all of the hidden patterns

Measure the Quality of Clustering

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

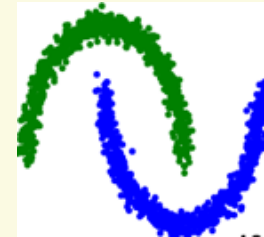
Considerations for Cluster Analysis

- ✓ Partitioning criteria
 - Single level vs. hierarchical partitioning (often, multi-level hierarchical partitioning is desirable)
- ✓ Separation of clusters
 - Exclusive (e.g., one point to only one region) vs. non-exclusive (e.g., one point to more than one class)
- ✓ Similarity measure
 - Distance-based (e.g., Euclidean, Manhattan, network, vector) vs. connectivity-based (e.g., graph, fuzzy, ambiguity)
- ✓ Clustering space
 - Full space (often when low dimensional) vs. subspaces (often in high-dimensional clustering)



Requirements and Challenges

- ✓ Scalability
 - Clustering all the data instead of only on samples
- ✓ Ability to deal with different types of attributes
 - Numerical, binary, categorical, ordinal, linked, and mixture of these
- ✓ Constraint-based clustering
 - User may give inputs on constraints
 - Use domain knowledge to determine input parameters
- ✓ Interpretability and usability
- ✓ Others
 - Discovery of clusters with arbitrary shape
 - Ability to deal with noisy data
 - Incremental clustering and insensitivity to input order
 - High dimensionality



Major Clustering Approaches (I)

- ✓ Partitioning approach:
 - Construct various partitions and then evaluate them by some criterion, e.g., minimizing the sum of square errors
 - Typical methods: k-means, k-medoids
- ✓ Hierarchical approach:
 - Create a hierarchical decomposition of the set of data (or objects) using some criterion (agglomerative or divisive)
 - Typical methods: Diana, Agnes, BIRCH, CAMELEON
- ✓ Density-based approach:
 - Based on connectivity and density functions
 - Typical methods: DBSCAN, OPTICS, DenClue
- ✓ Grid-based approach:
 - based on a multiple-level granularity structure
 - Typical methods: STING, WaveCluster, CLIQUE

Major Clustering Approaches (II)

✓ Model-based:

- A model is hypothesized for each of the clusters and tries to find the best fit of that model to each other
- Typical methods: EM, SOM, COBWEB

✓ Frequent pattern-based:

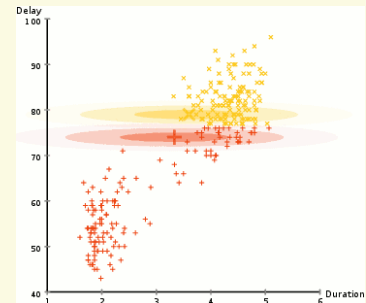
- Based on the analysis of frequent patterns
- Typical methods: p-Cluster

✓ User-guided or constraint-based:

- Clustering by considering user-specified or application-specific constraints
- Typical methods: COD (obstacles), constrained clustering

✓ Link-based clustering:

- Objects are often linked together in various ways
- Massive links can be used to cluster objects: SimRank, LinkClus



Data clustering – Partition-based

Partitioning Algorithms: Basic Concept

- ✓ Partitioning method: Partitioning a database ***D*** of ***n*** objects into a set of ***k*** clusters, such that the sum of squared distances is minimized (where c_i is the centroid or medoid of cluster C_i)

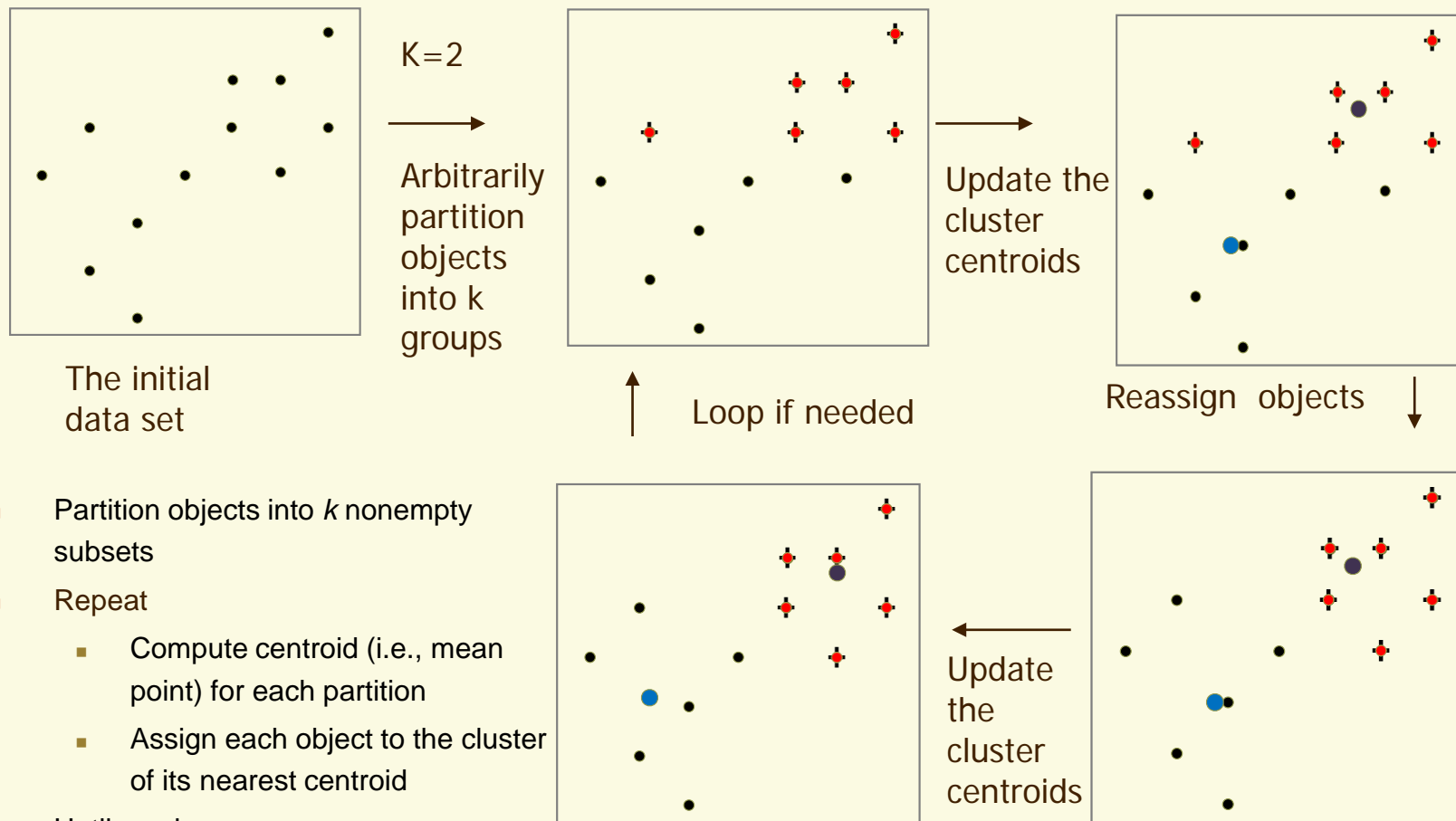
$$E = \sum_{i=1}^k \sum_{p \in C_i} (d(p, c_i))^2$$

- ✓ Given k , find a partition of k clusters that optimizes the chosen partitioning criterion
 - Global optimal: exhaustively enumerate all partitions
 - Heuristic methods: *k-means* and *k-medoids* algorithms
 - *k-means* (MacQueen'67, Lloyd'57/'82): Each cluster is represented by the center of the cluster
 - *k-medoids* or PAM (Partition around medoids) (Kaufman & Rousseeuw'87): Each cluster is represented by one of the objects in the cluster

The *K-Means* Clustering Method

- ✓ Given k , the *k-means* algorithm is implemented in four steps:
- Partition objects into k nonempty subsets
 - Compute seed points as the centroids of the clusters of the current partitioning (the centroid is the center, i.e., *mean point*, of the cluster)
 - Assign each object to the cluster with the nearest seed point
 - Go back to Step 2, stop when the assignment does not change

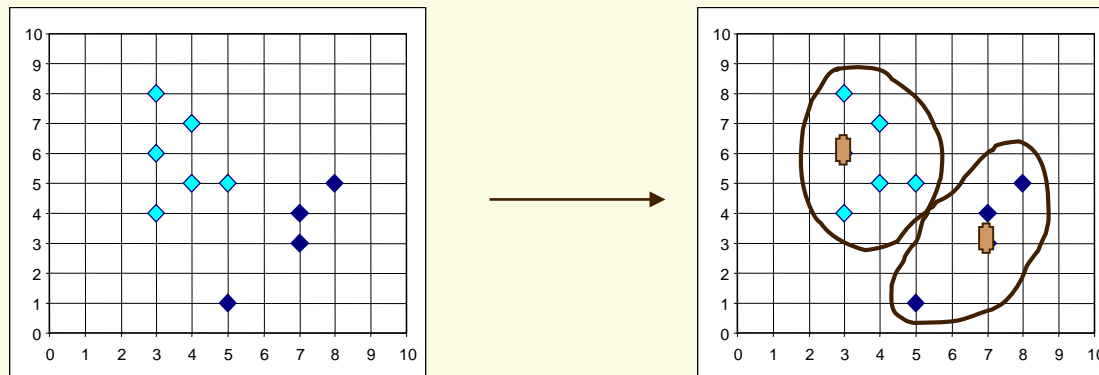
K-Means Clustering



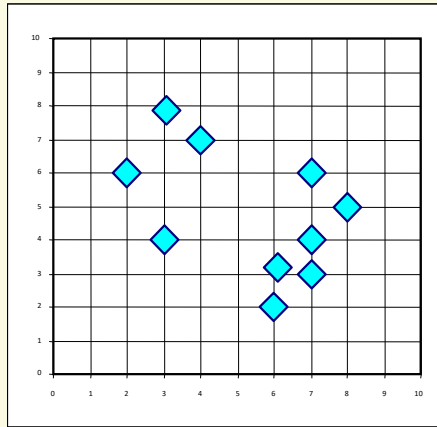
- Partition objects into k nonempty subsets
- Repeat
 - Compute centroid (i.e., mean point) for each partition
 - Assign each object to the cluster of its nearest centroid
- Until no change

What Is the Problem of the K-Means Method?

- ✓ The k-means algorithm is sensitive to outliers !
 - Since an object with an extremely large value may substantially distort the distribution of the data
- ✓ K-Medoids: Instead of taking the **mean** value of the object in a cluster as a reference point, **medoids** can be used, which is the **most centrally located** object in a cluster

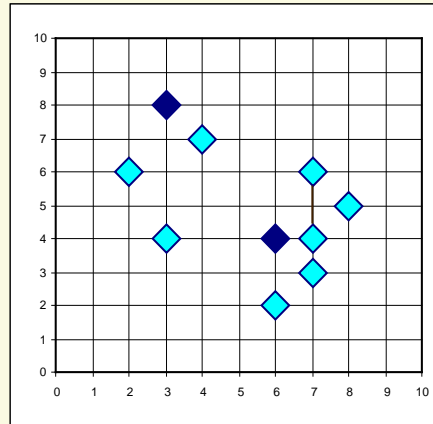


A Typical K-Medoids Algorithm



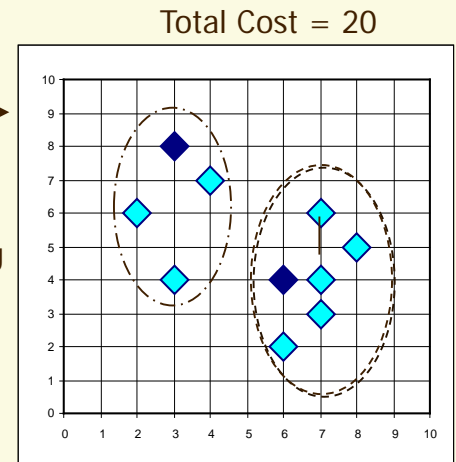
K=2

Arbitrary
choose k
object as
initial
medoids

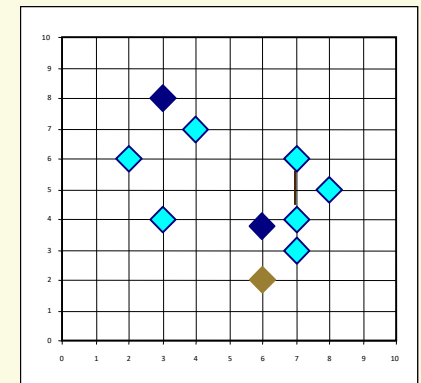


Total Cost = 26

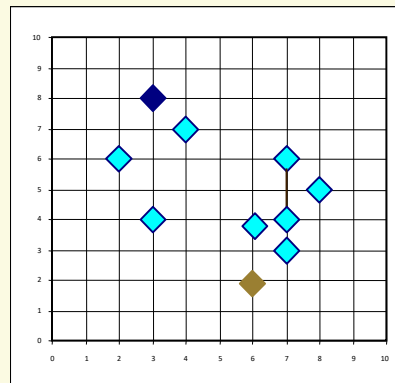
Assign
each
remaining
object to
nearest
medoids



Randomly select a
nonmedoid object, O_{random}



Compute
total cost of
swapping



$$E = \sum_{i=1}^k \sum_{p \in C_i} (d(p, o_i))$$

**Do loop
Until no
change**

Swapping O
and O_{random}
If quality is
improved.

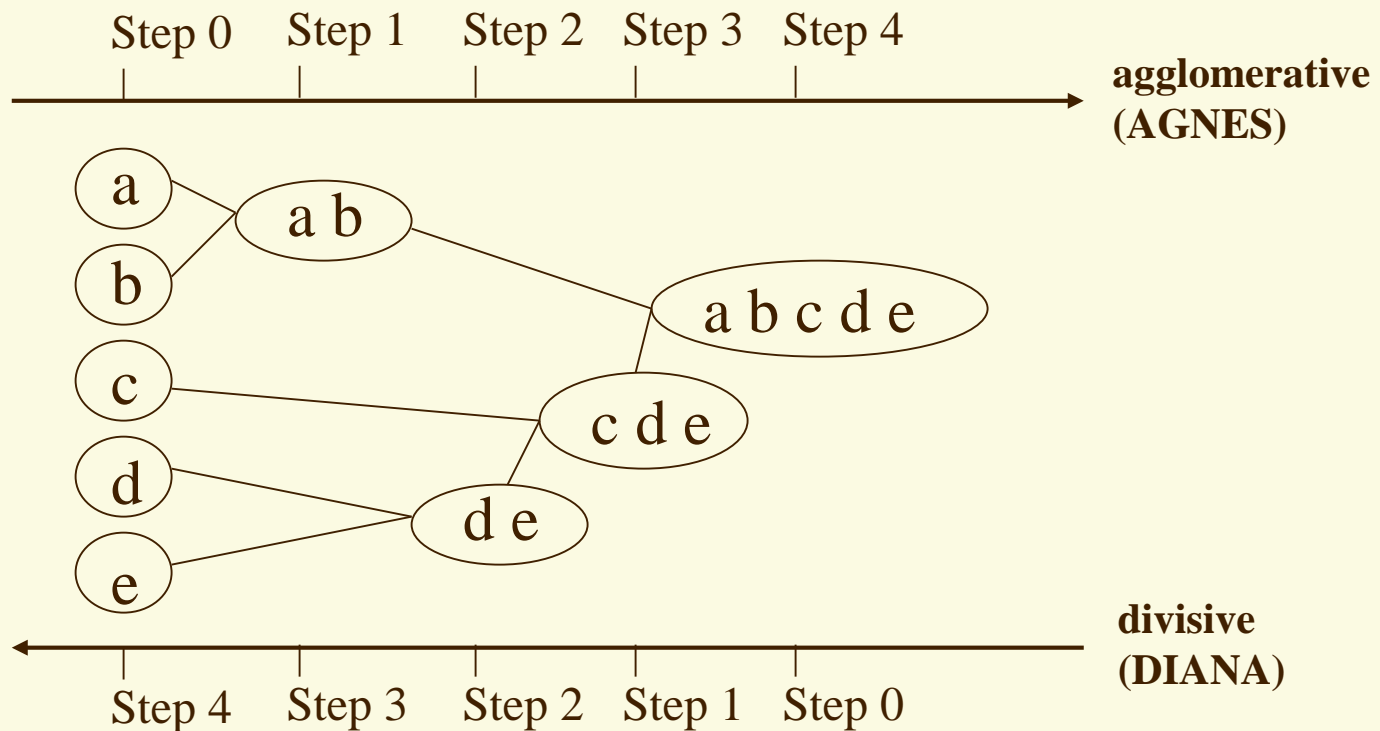
The K-Medoid Clustering Method

- ✓ *K-Medoids* Clustering: Find *representative* objects (medoids) in clusters
 - *PAM* (Partitioning Around Medoids, Kaufmann & Rousseeuw 1987)
 - Starts from an initial set of medoids and iteratively replaces one of the medoids by one of the non-medoids if it improves the total distance of the resulting clustering
 - *PAM* works effectively for small data sets, but does not scale well for large data sets (due to the computational complexity)
- ✓ Efficiency improvement on PAM
 - *CLARA* (Kaufmann & Rousseeuw, 1990): PAM on samples
 - *CLARANS* (Ng & Han, 1994): Randomized re-sampling

Data clustering – Hierarchical methods

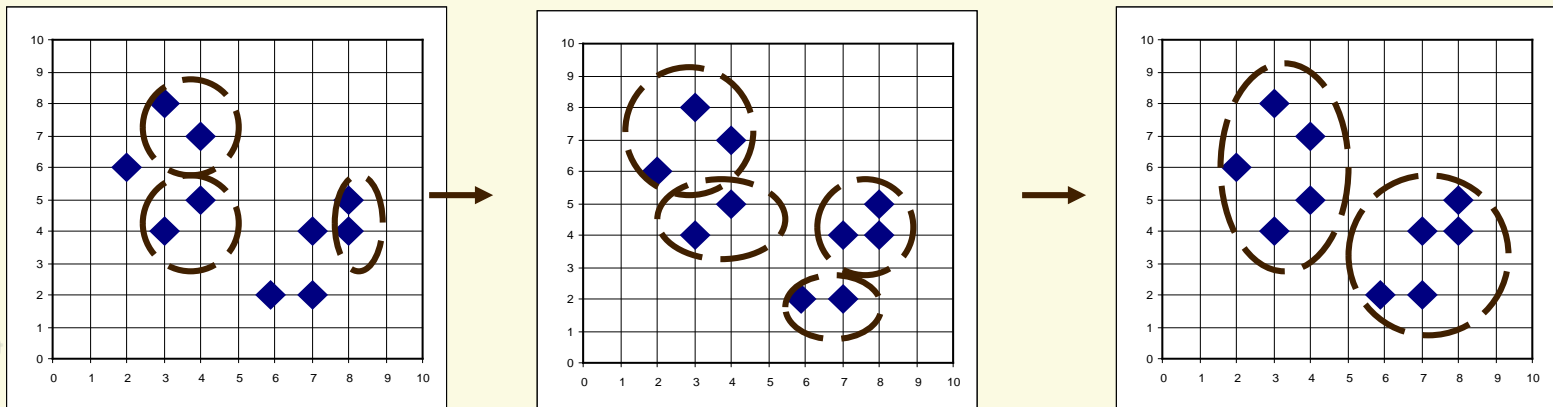
Hierarchical Clustering

- ✓ Use distance matrix as clustering criteria. This method does not require the number of clusters k as an input, but needs a termination condition



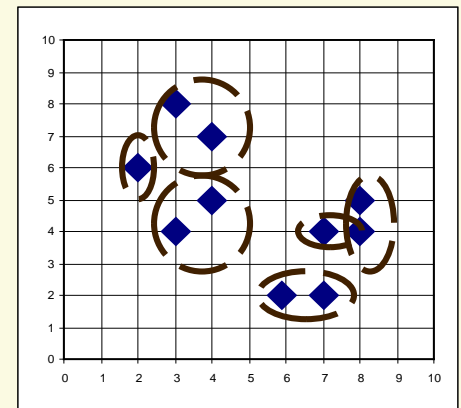
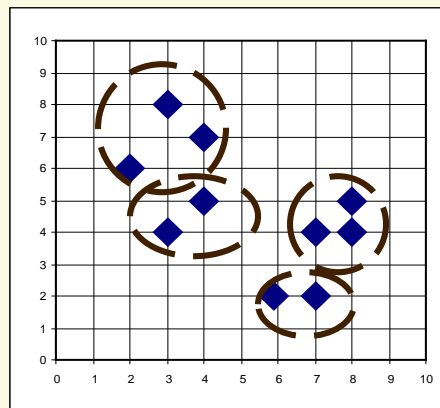
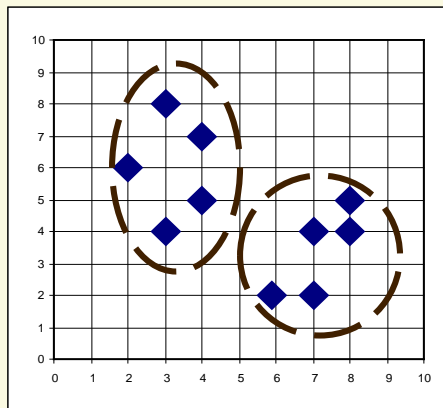
AGNES (Agglomerative Nesting)

- ✓ Introduced in Kaufmann and Rousseeuw (1990)
- ✓ Implemented in statistical packages, e.g., Splus
- ✓ Use the **single-link** method and the dissimilarity matrix
- ✓ Merge nodes that have the least dissimilarity
- ✓ Go on in a non-descending fashion
- ✓ Eventually all nodes belong to the same cluster



DIANA (Divisive Analysis)

- ✓ Introduced in Kaufmann and Rousseeuw (1990)
- ✓ Implemented in statistical analysis packages, e.g., Splus
- ✓ Inverse order of AGNES
- ✓ Eventually each node forms a cluster on its own



Data clustering – Density-based methods

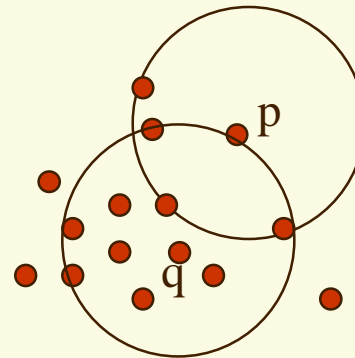
Density-Based Clustering Methods

- ✓ Clustering based on density (local cluster criterion), such as density-connected points
- ✓ Major features:
 - Discover clusters of arbitrary shape
 - Handle noise
 - One scan
 - Need density parameters as termination condition
- ✓ Several interesting studies:
 - DBSCAN: Ester, et al. (KDD'96)
 - OPTICS: Ankerst, et al (SIGMOD'99).
 - DENCLUE: Hinneburg & D. Keim (KDD'98)
 - CLIQUE: Agrawal, et al. (SIGMOD'98) (more grid-based)

Density-Based Clustering: Basic Concepts

- ✓ Two parameters:
 - *Eps*: Maximum radius of the neighborhood
 - *MinPts*: Minimum number of points in an Eps-neighborhood of that point
- ✓ $N_{Eps}(q)$: $\{p \text{ belongs to } D \mid \text{dist}(p,q) \leq Eps\}$
- ✓ **Directly density-reachable**: A point p is directly density-reachable from a point q w.r.t. Eps , $MinPts$ if
 - p belongs to $N_{Eps}(q)$
 - q is a core point:

$$|N_{Eps}(q)| \geq MinPts$$



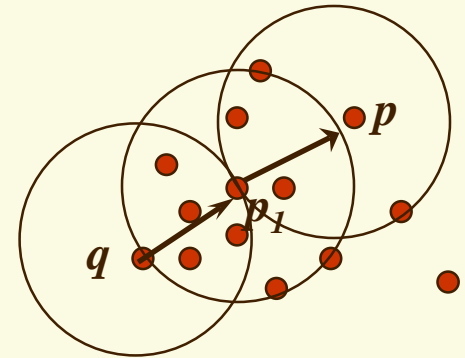
$MinPts = 5$

$Eps = 1 \text{ cm}$

Density-Reachable and Density-Connected

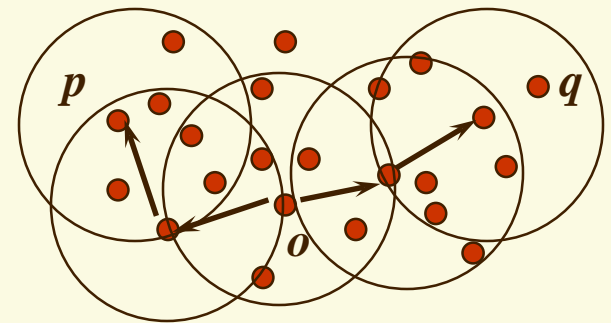
✓ Density-reachable:

- A point p is **density-reachable** from a point q w.r.t. Eps , $MinPts$ if there is a chain of points p_1, \dots, p_n , $p_1 = q$, $p_n = p$ such that p_{i+1} is directly density-reachable from p_i



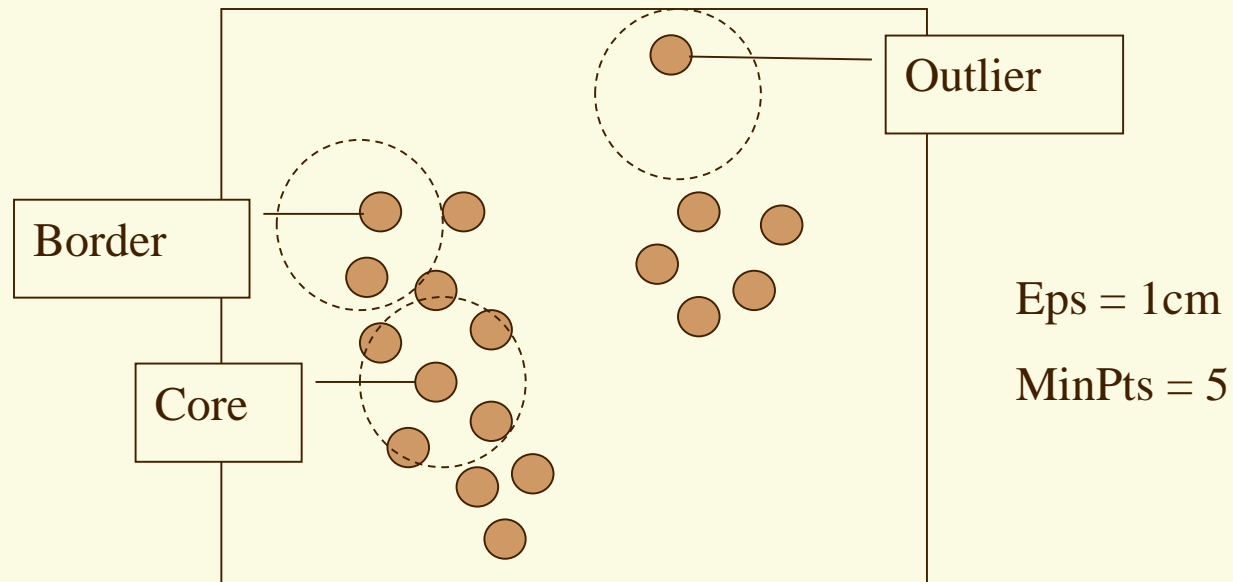
✓ Density-connected

- A point p is **density-connected** to a point q w.r.t. Eps , $MinPts$ if there is a point o such that both, p and q are density-reachable from o w.r.t. Eps and $MinPts$



DBSCAN: Density-Based Spatial Clustering of Applications with Noise

- ✓ Relies on a *density-based* notion of cluster: A *cluster* is defined as a maximal set of density-connected points
- ✓ Discovers clusters of arbitrary shape in spatial databases with noise



DBSCAN: Sensitive to Parameters

Figure 8. DBScan results for DS1 with MinPts at 4 and Eps at (a) 0.5 and (b) 0.4.

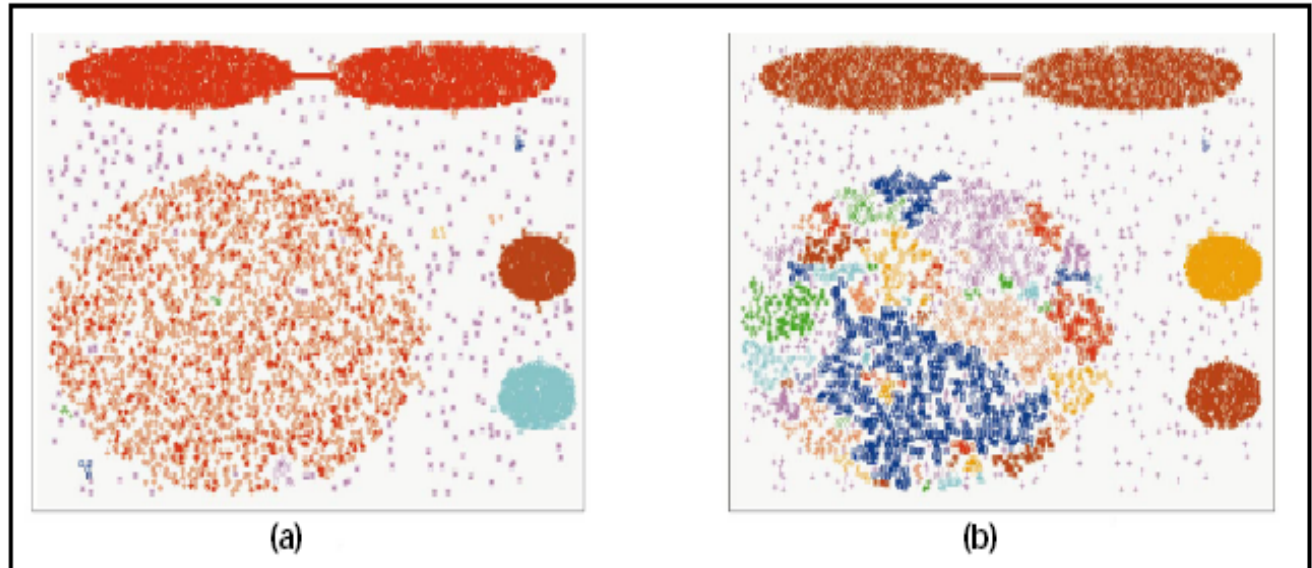
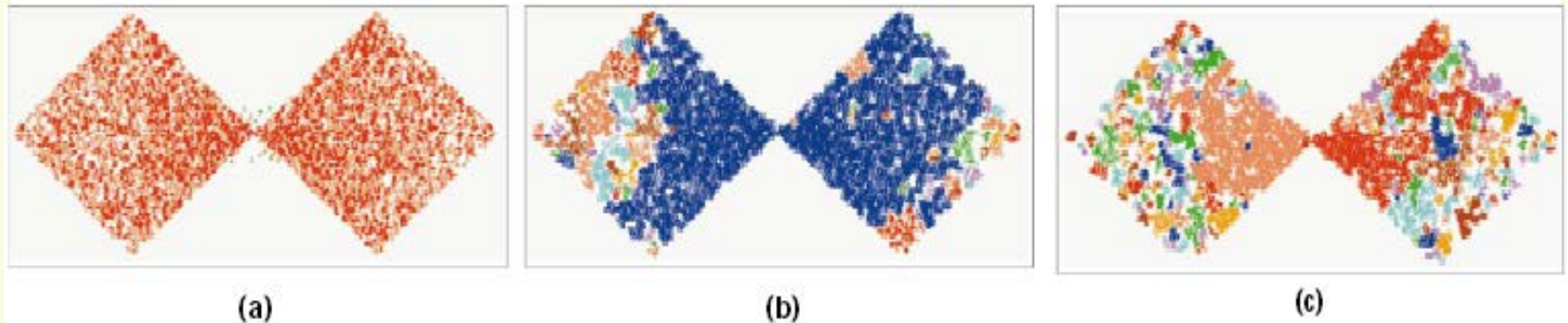


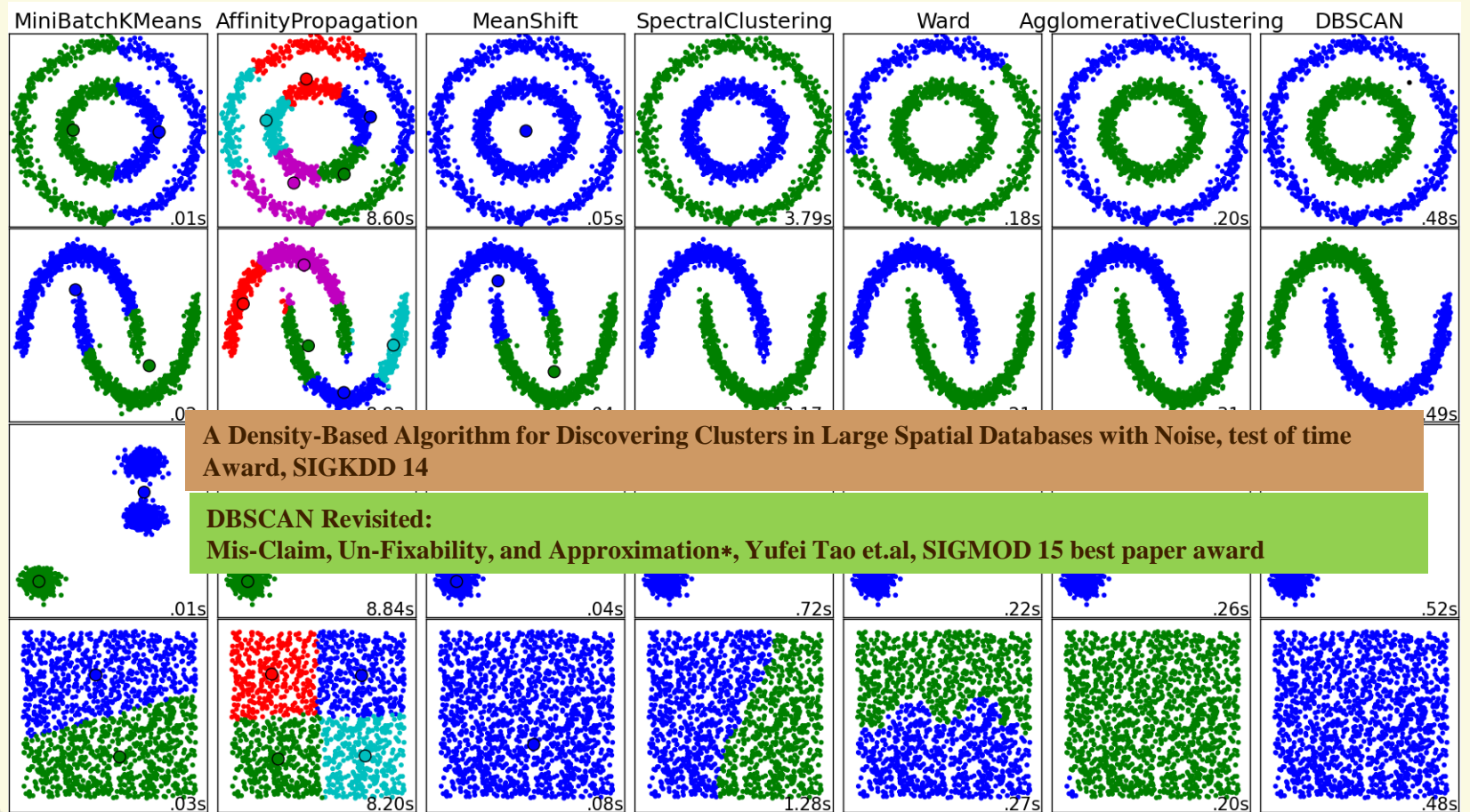
Figure 9. DBScan results for DS2 with MinPts at 4 and Eps at (a) 5.0, (b) 3.5, and (c) 3.0.



DBSCAN online Demo:

<http://webdocs.cs.ualberta.ca/~yaling/Cluster/Applet/Code/Cluster.html>

Comparing different clustering algorithms



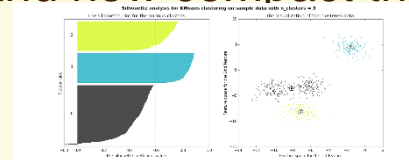
Data clustering – evaluation

Determine the Number of Clusters

- ✓ Empirical method
 - # of clusters: $k \approx \sqrt{n/2}$ for a dataset of n points, e.g., $n = 200$, $k = 10$
- ✓ Elbow method
 - Use the turning point in the curve of sum of within cluster variance w.r.t the # of clusters
- ✓ Cross validation method
 - Divide a given data set into m parts
 - Use $m - 1$ parts to obtain a clustering model
 - Use the remaining part to test the quality of the clustering
 - E.g., For each point in the test set, find the closest centroid, and use the sum of squared distance between all points in the test set and the closest centroids to measure how well the model fits the test set
 - For any $k > 0$, repeat it m times, compare the overall quality measure w.r.t. different k 's, and find # of clusters that fits the data the best

Measuring Clustering Quality

- ✓ How do I know whether the clustering results are good?
- ✓ 3 kinds of measures: External, internal and relative
- ✓ External: supervised, employ criteria not inherent to the dataset
 - Compare a clustering against prior or expert-specified knowledge (i.e., the ground truth) using certain clustering quality measure
- ✓ Internal: unsupervised, criteria derived from data itself
 - Evaluate the goodness of a clustering by considering how well the clusters are separated, and how compact the clusters are, e.g., **Silhouette coefficient**
- ✓ Relative: directly compare different clusterings, usually those obtained via different parameter settings for the same algorithm

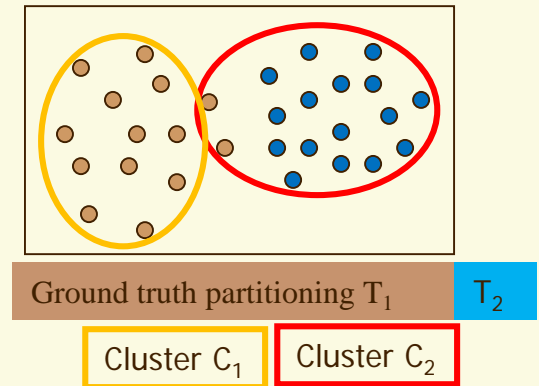


Measuring Clustering Quality: External Methods

- ✓ Clustering quality measure: $Q(C, T)$, for a clustering C given the ground truth T
- ✓ Q is good if it satisfies the following 4 essential criteria
 - Cluster homogeneity: the purer, the better
 - Cluster completeness: should assign objects belong to the same category in the ground truth to the same cluster
 - Rag bag: putting a heterogeneous object into a pure cluster should be penalized more than putting it into a *rag bag* (i.e., “miscellaneous” or “other” category)
 - Small cluster preservation: splitting a small category into pieces is more harmful than splitting a large category into pieces

Some Commonly Used External Measures

- ✓ Matching-based measures
 - Purity, maximum matching, F-measure
- ✓ Entropy-Based Measures
 - Conditional entropy, normalized mutual information (NMI), variation of information
- ✓ Pair-wise measures
 - Four possibilities: True positive (TP), FN, FP, TN
 - Jaccard coefficient, Rand statistic, Fowlkes-Mallow measure
- ✓ Correlation measures
 - Discretized Huber static, normalized discretized Huber static



Graph Clustering

Clustering Graphs and Network Data

✓ Applications

- Bi-partite graphs, e.g., customers and products, authors and conferences
- Web search engines, e.g., click through graphs and Web graphs
- Social networks, friendship/coauthor graphs

✓ Similarity measures

- Geodesic distances
- Distance based on random walk (SimRank)

✓ Graph clustering methods

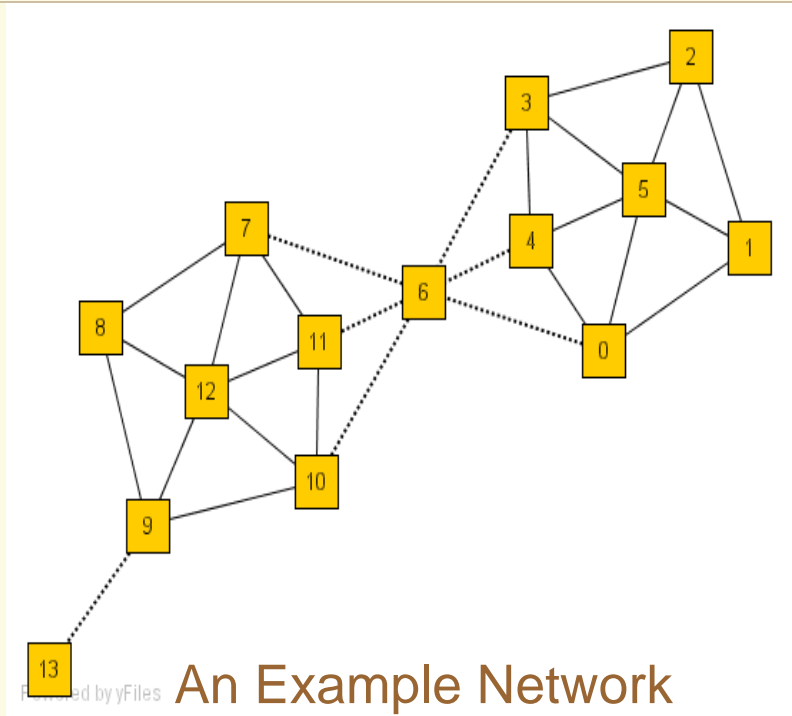
- Minimum cuts: FastModularity (Clauset, Newman & Moore, 2004)
- Density-based clustering: SCAN (Xu et al., KDD'2007)

Two Approaches for Graph Clustering

- ✓ Two approaches for clustering graph data
 - Use *generic clustering methods* for high-dimensional data
 - *Designed specifically for clustering graphs*
- ✓ Using clustering methods for high-dimensional data
 - Extract a similarity matrix from a graph using a similarity measure
 - A generic clustering method can then be applied on the similarity matrix to discover clusters
 - Spectral clustering methods: approximate optimal graph cuts
- ✓ Methods specific to graphs
 - Search the graph to find well-connected components as clusters
 - Ex. SCAN (Structural Clustering Algorithm for Networks)
 - X. Xu, N. Yuruk, Z. Feng, and T. A. J. Schweiger, “SCAN: A Structural Clustering Algorithm for Networks”, KDD'07

SCAN: Density-Based Clustering of Networks

- ✓ How many clusters?
- ✓ What size should they be?
- ✓ What is the best partitioning?
- ✓ Should some points be segregated?



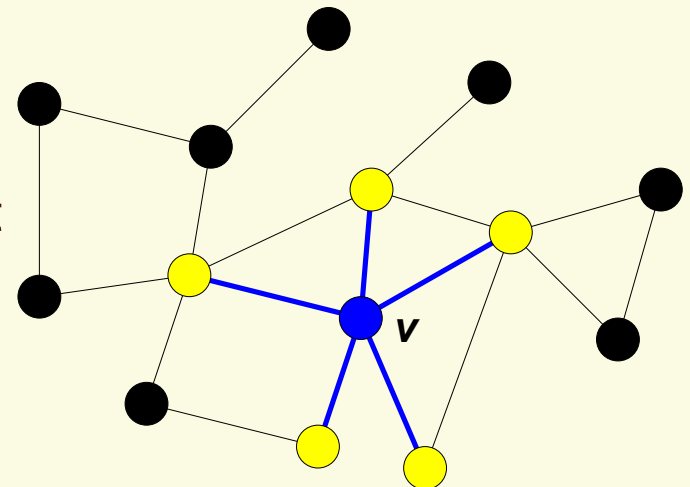
- Application: Given simply information of who associates with whom, could one identify clusters of individuals with common interests or special relationships (families, cliques, terrorist cells)?

A Social Network Model

- ✓ Cliques, hubs and outliers
 - Individuals in a tight social group, or **clique**, know many of the same people, regardless of the size of the group
 - Individuals who are **hubs** know many people in different groups but belong to no single group. Politicians, for example bridge multiple groups
 - Individuals who are **outliers** reside at the margins of society. Hermits, for example, know few people and belong to no group

- ✓ The Neighborhood of a Vertex

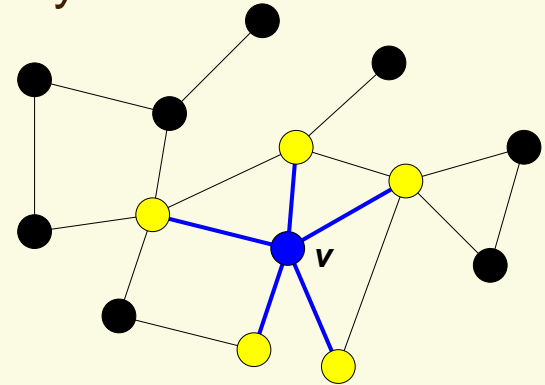
- Define $\Gamma(v)$ as the immediate neighborhood of a vertex (i.e. the set of people that an individual knows)



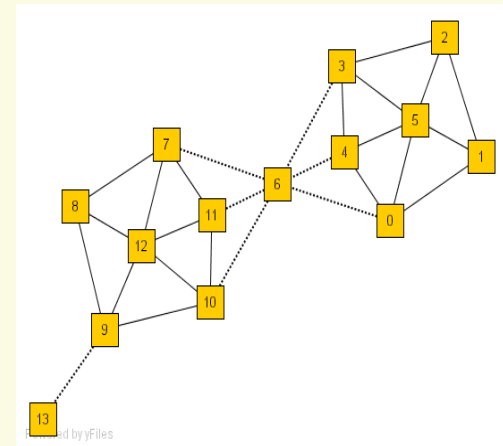
Structure Similarity

- ✓ The desired features tend to be captured by a measure we call Structural Similarity

$$\sigma(v, w) = \frac{|\Gamma(v) \cap \Gamma(w)|}{\sqrt{|\Gamma(v)| |\Gamma(w)|}}$$



- ✓ Structural similarity is large for members of a clique and small for hubs and outliers



Structural Connectivity [1]

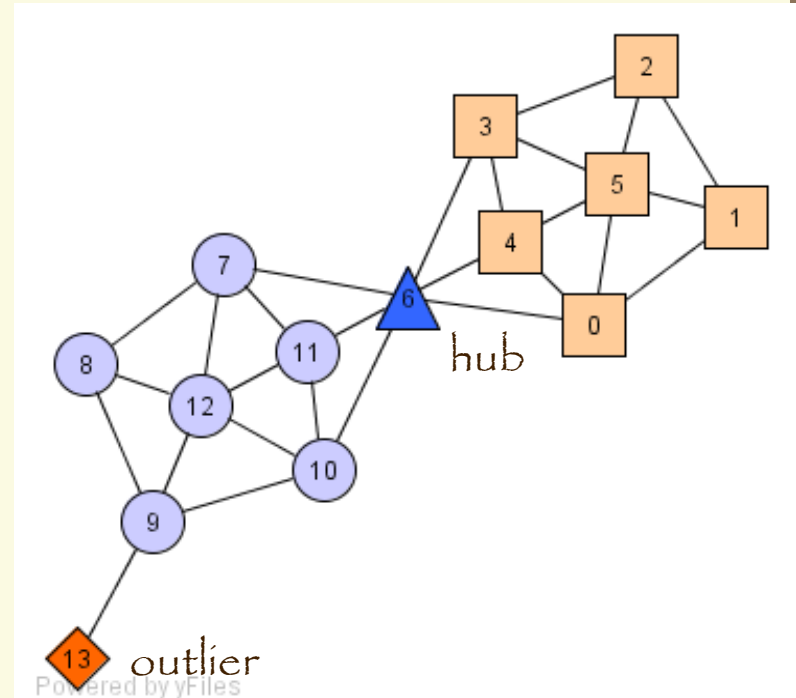
- ✓ ε -Neighborhood: $N_\varepsilon(v) = \{w \in \Gamma(v) \mid \sigma(v, w) \geq \varepsilon\}$
- ✓ Core: $CORE_{\varepsilon, \mu}(v) \Leftrightarrow |N_\varepsilon(v)| \geq \mu$
- ✓ Direct structure reachable:
$$DirRECH_{\varepsilon, \mu}(v, w) \Leftrightarrow CORE_{\varepsilon, \mu}(v) \wedge w \in N_\varepsilon(v)$$
- ✓ Structure reachable: transitive closure of direct structure reachability
- ✓ Structure connected:

$$CONNECT_{\varepsilon, \mu}(v, w) \Leftrightarrow \exists u \in V : RECH_{\varepsilon, \mu}(u, v) \wedge RECH_{\varepsilon, \mu}(u, w)$$

[1] M. Ester, H. P. Kriegel, J. Sander, & X. Xu (KDD'96) "A Density-Based Algorithm for Discovering Clusters in Large Spatial Databases"

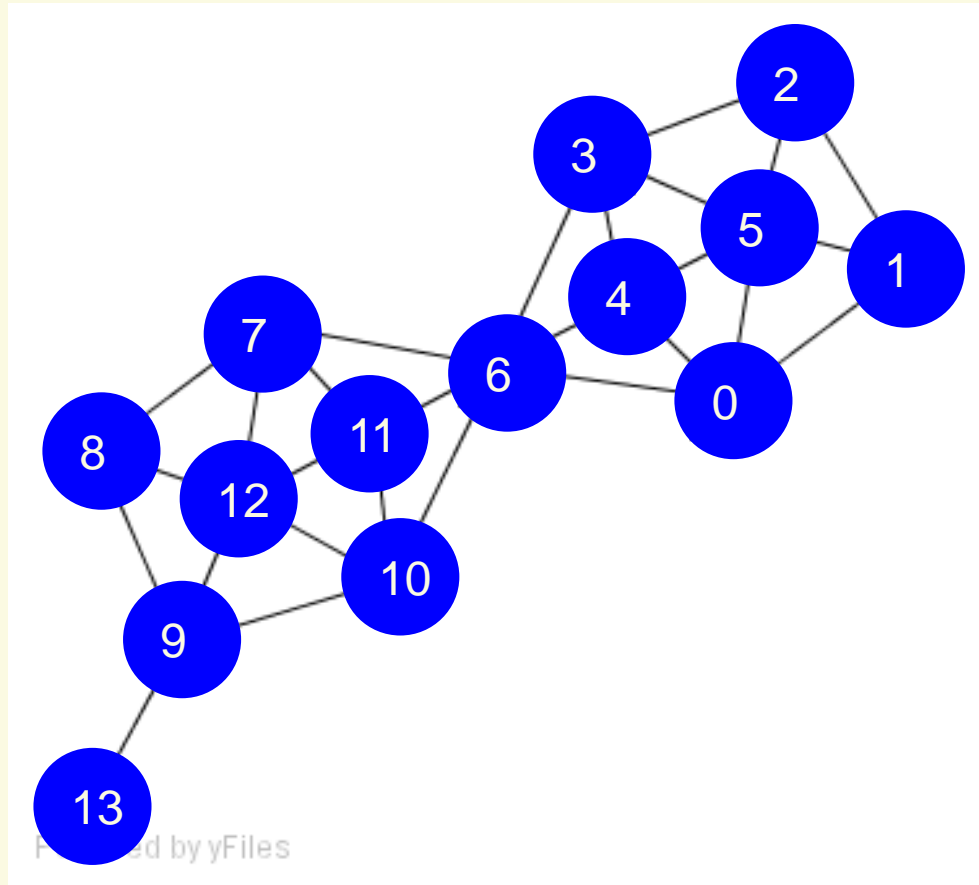
Structure-Connected Clusters

- ✓ Structure-connected cluster C
 - Connectivity: $\forall v, w \in C : CONNECT_{\varepsilon, \mu}(v, w)$
 - Maximality: $\forall v, w \in V : v \in C \wedge REACH_{\varepsilon, \mu}(v, w) \Rightarrow w \in C$
- ✓ Hubs:
 - Not belong to any cluster
 - Bridge to many clusters
- ✓ Outliers:
 - Not belong to any cluster
 - Connect to less clusters



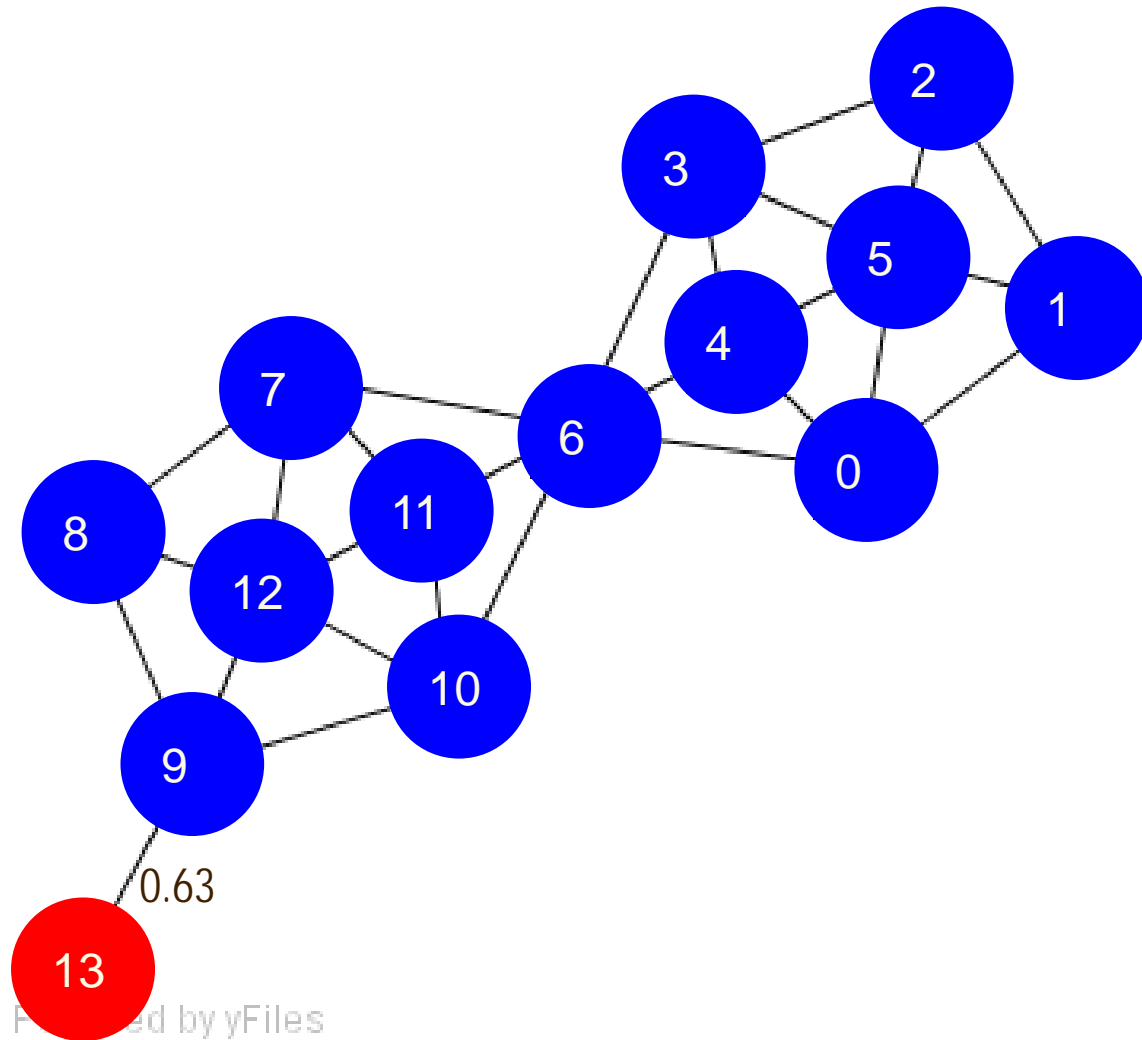
Algorithm

$$\mu = 2$$
$$\varepsilon = 0.7$$



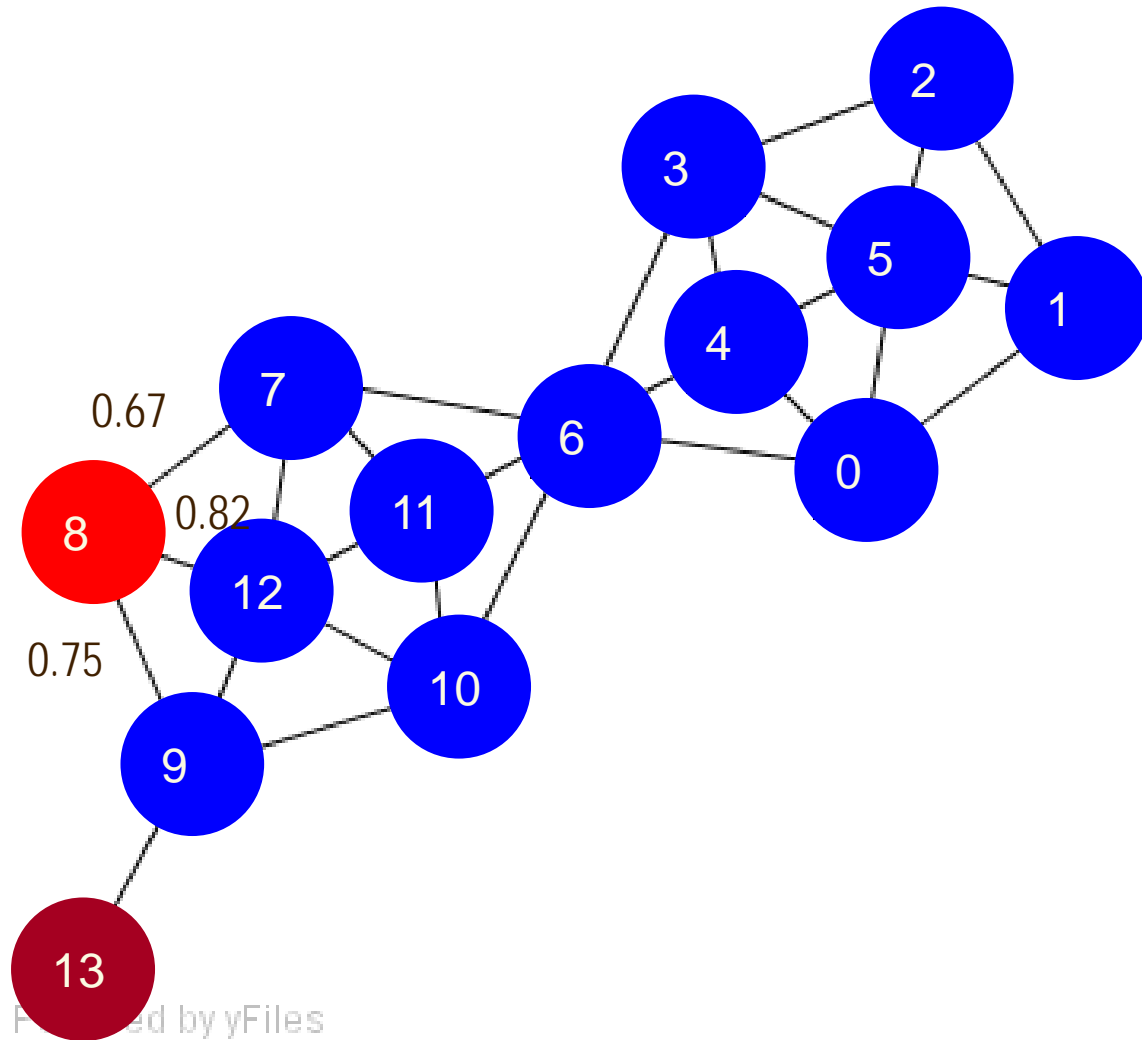
Algorithm

$$\mu = 2$$
$$\varepsilon = 0.7$$



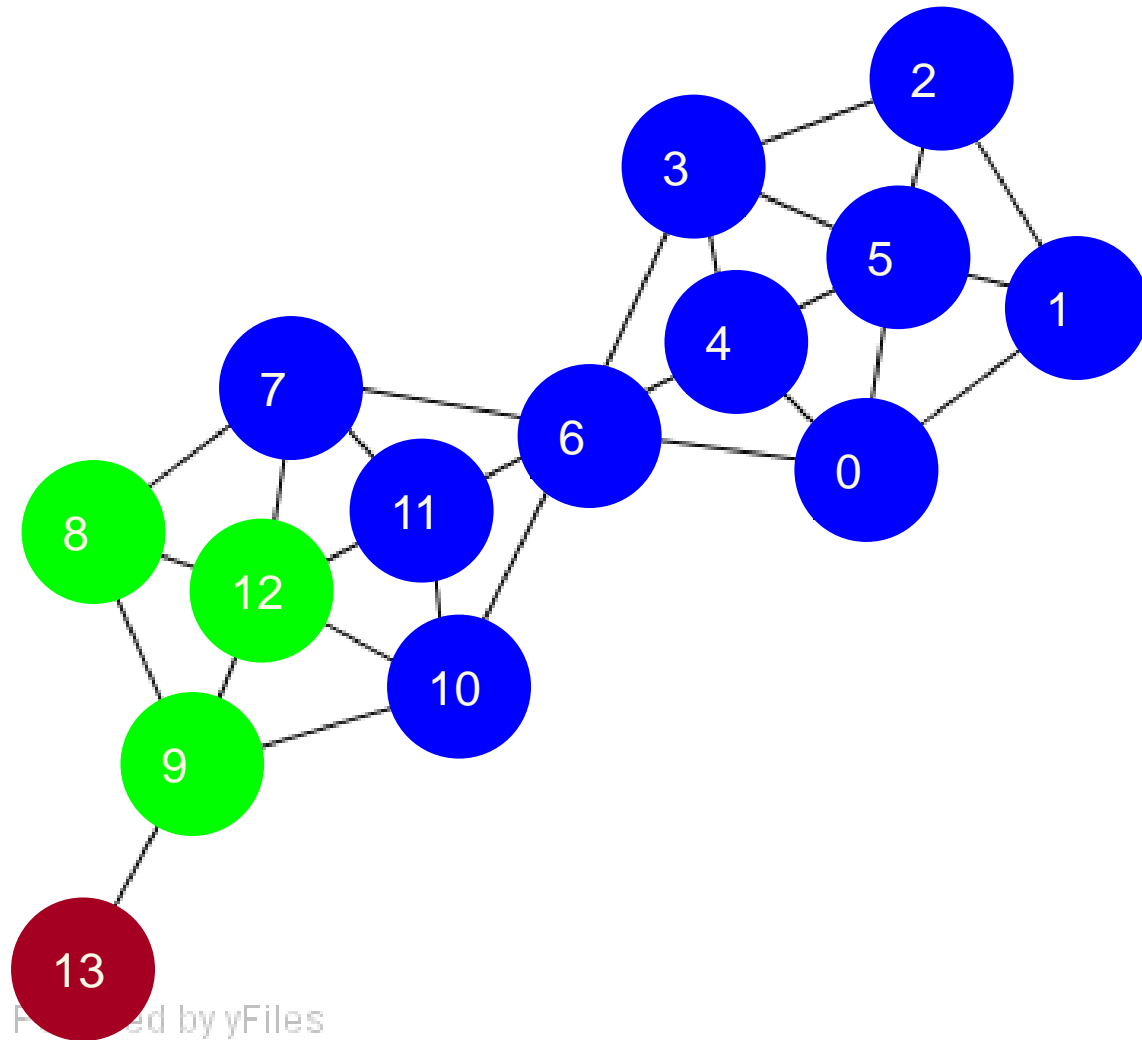
Algorithm

$$\mu = 2$$
$$\varepsilon = 0.7$$



Algorithm

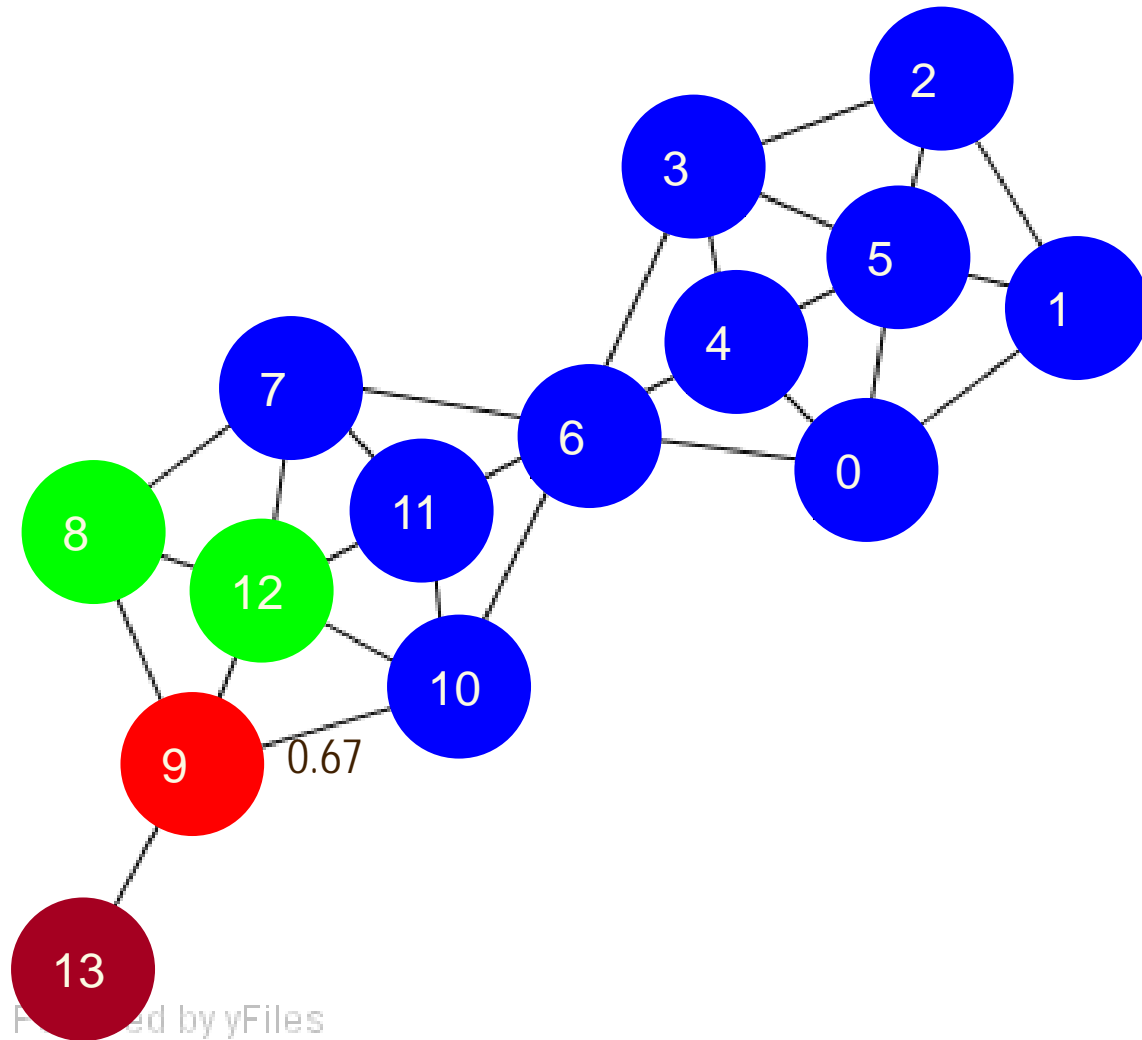
$$\mu = 2$$
$$\varepsilon = 0.7$$



Powered by yFiles

Algorithm

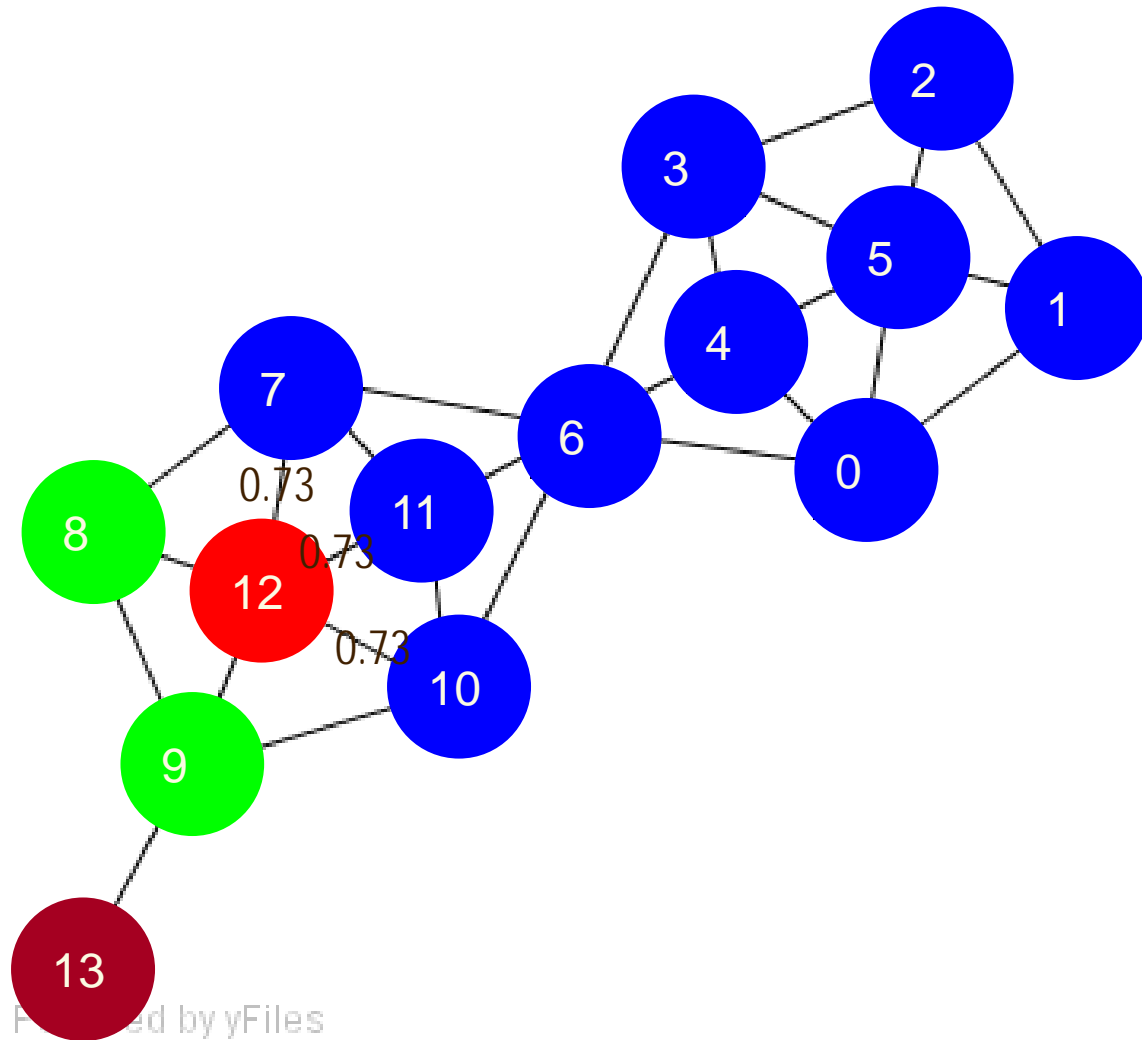
$$\mu = 2$$
$$\varepsilon = 0.7$$



Powered by yFiles

Algorithm

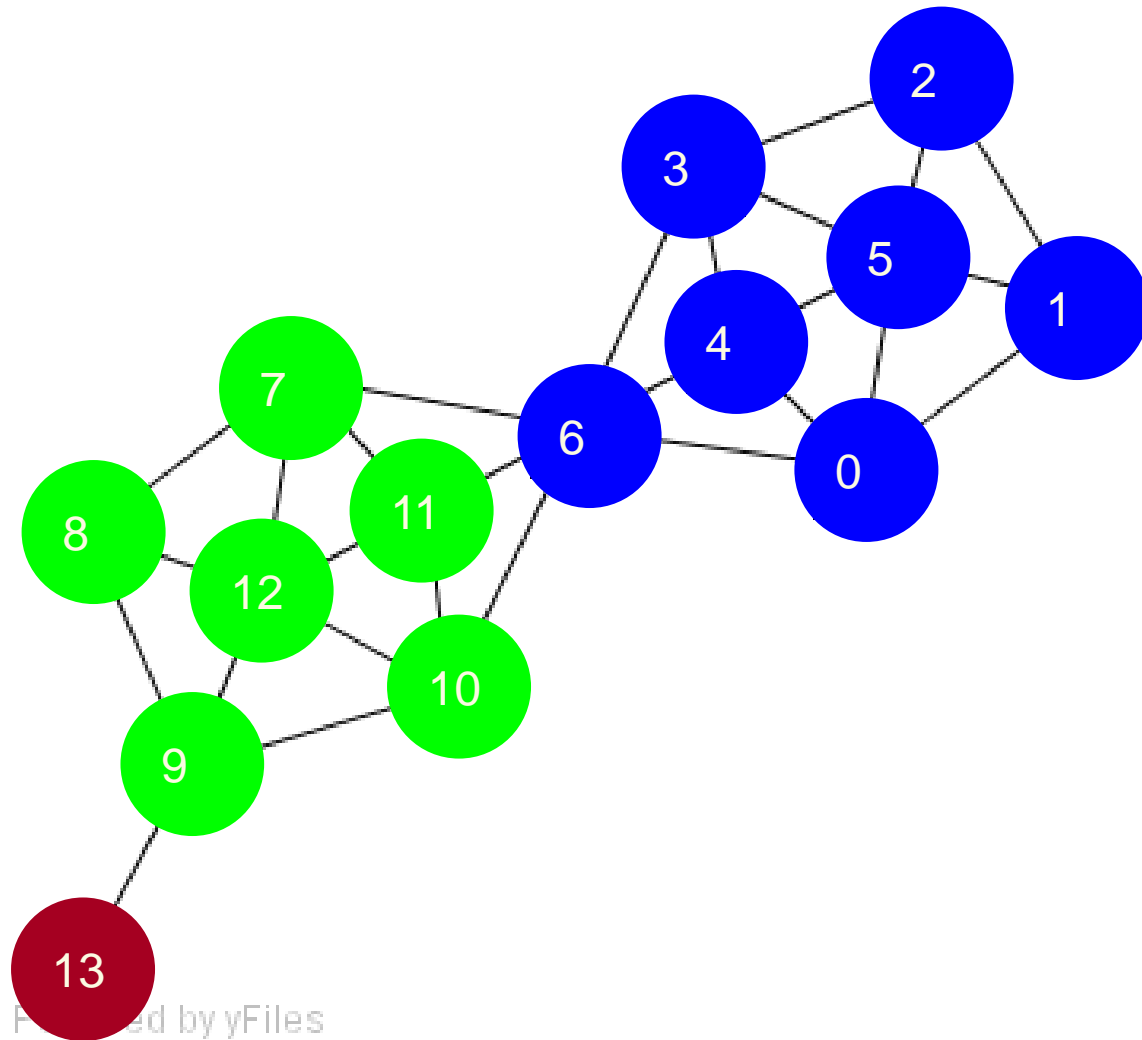
$$\mu = 2$$
$$\varepsilon = 0.7$$



Powered by yFiles

Algorithm

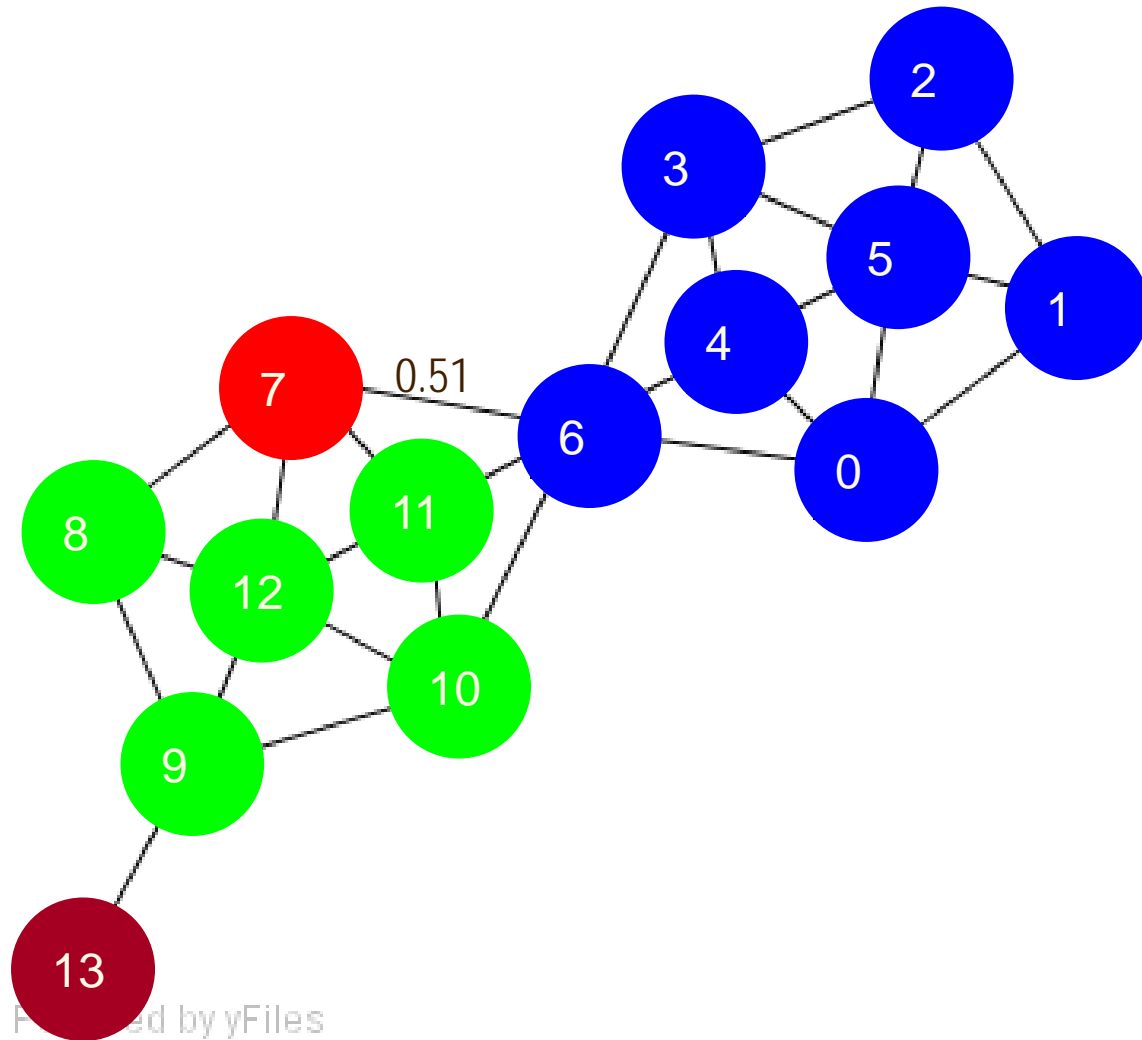
$$\mu = 2$$
$$\varepsilon = 0.7$$



Powered by yFiles

Algorithm

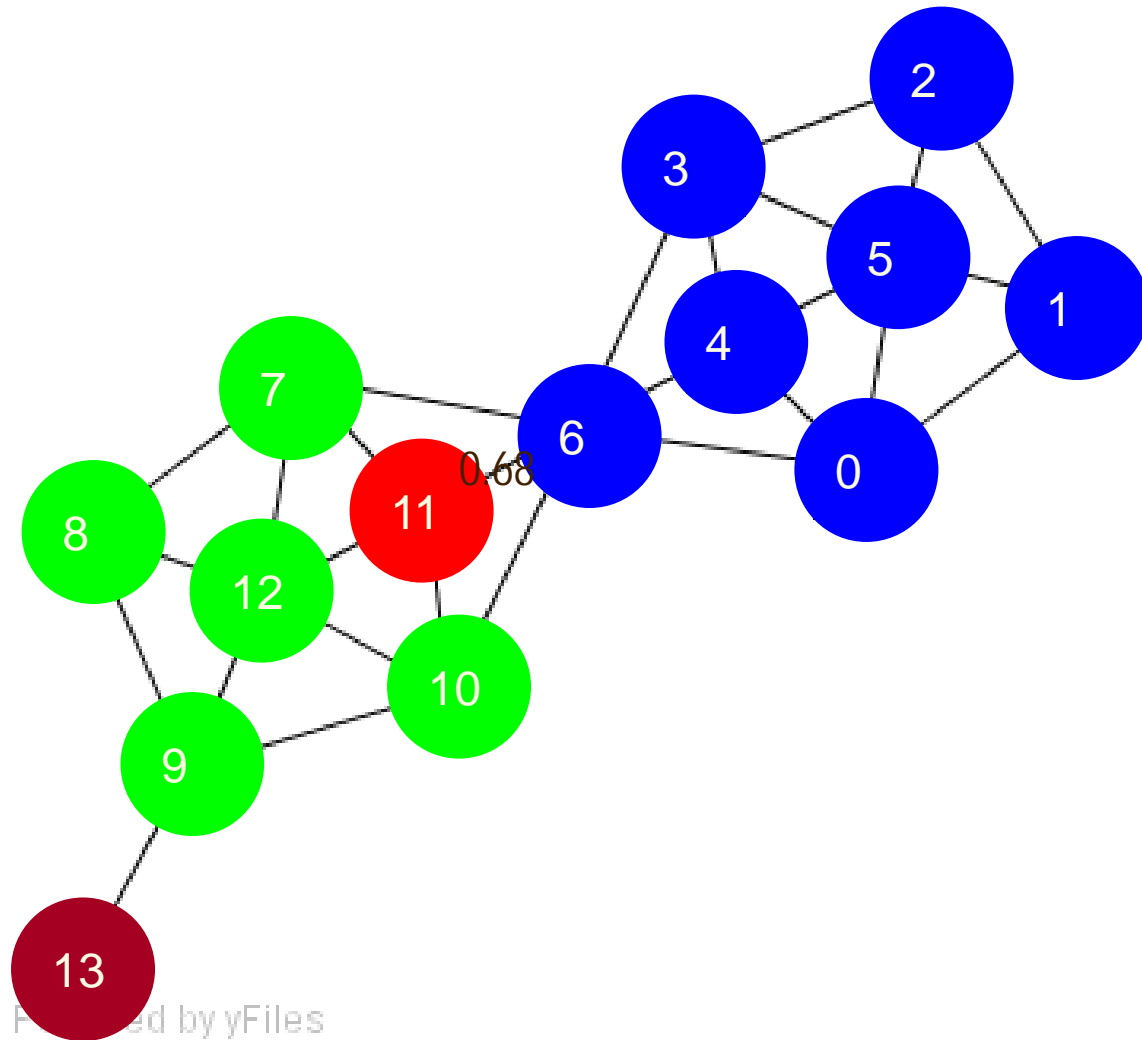
$$\mu = 2$$
$$\varepsilon = 0.7$$



Powered by yFiles

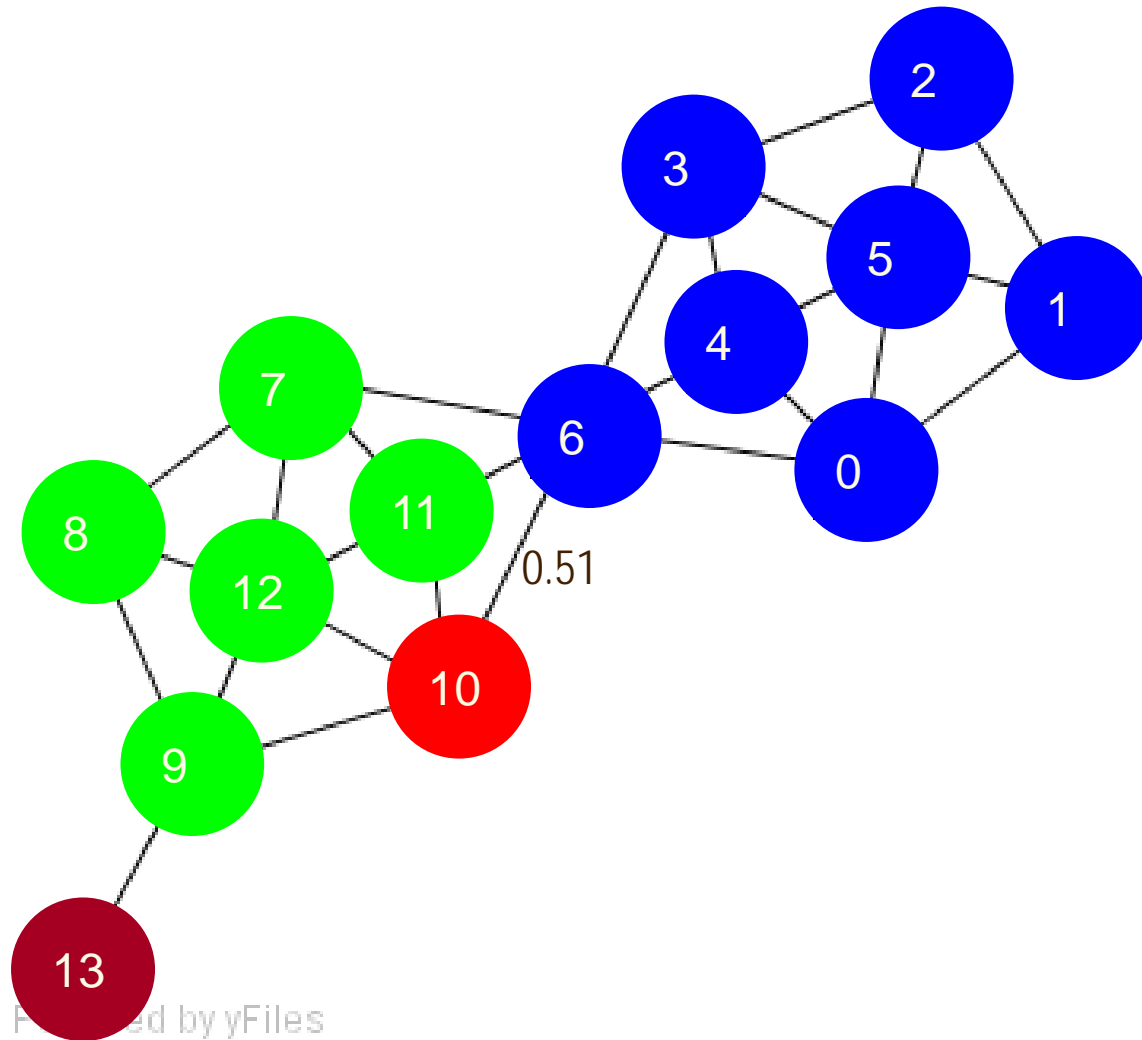
Algorithm

$$\mu = 2$$
$$\varepsilon = 0.7$$



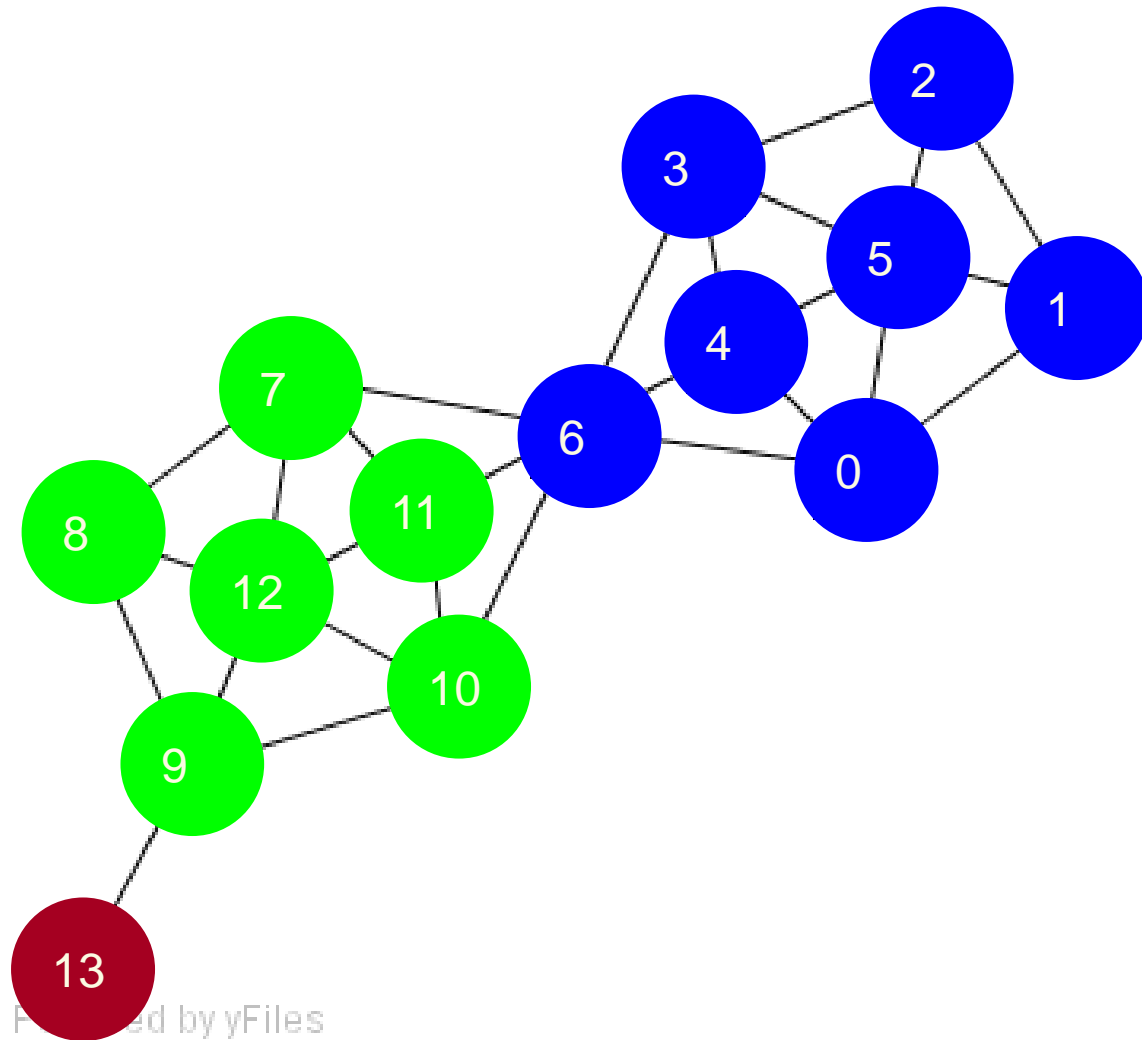
Algorithm

$$\mu = 2$$
$$\varepsilon = 0.7$$



Algorithm

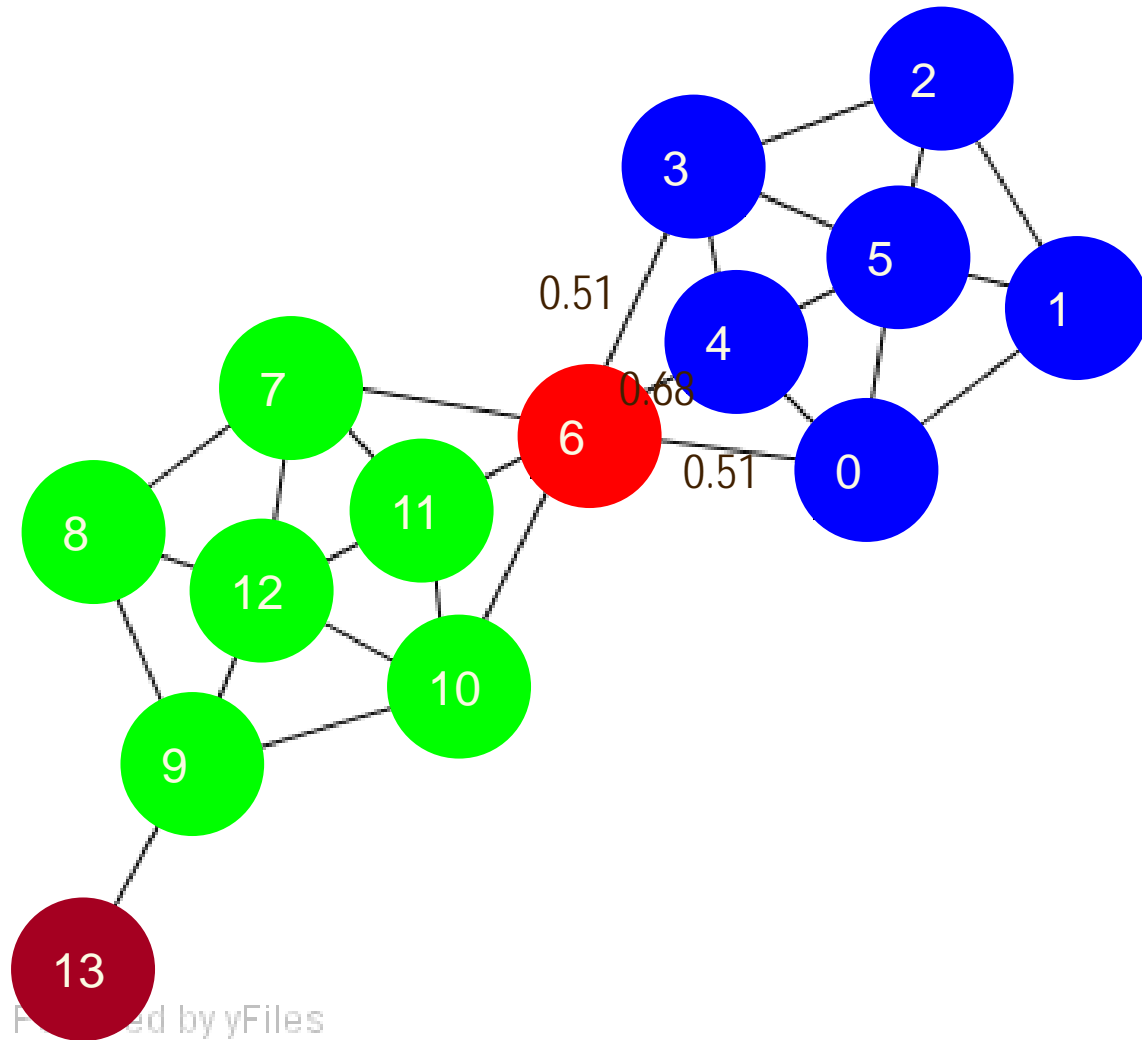
$$\mu = 2$$
$$\varepsilon = 0.7$$



Powered by yFiles

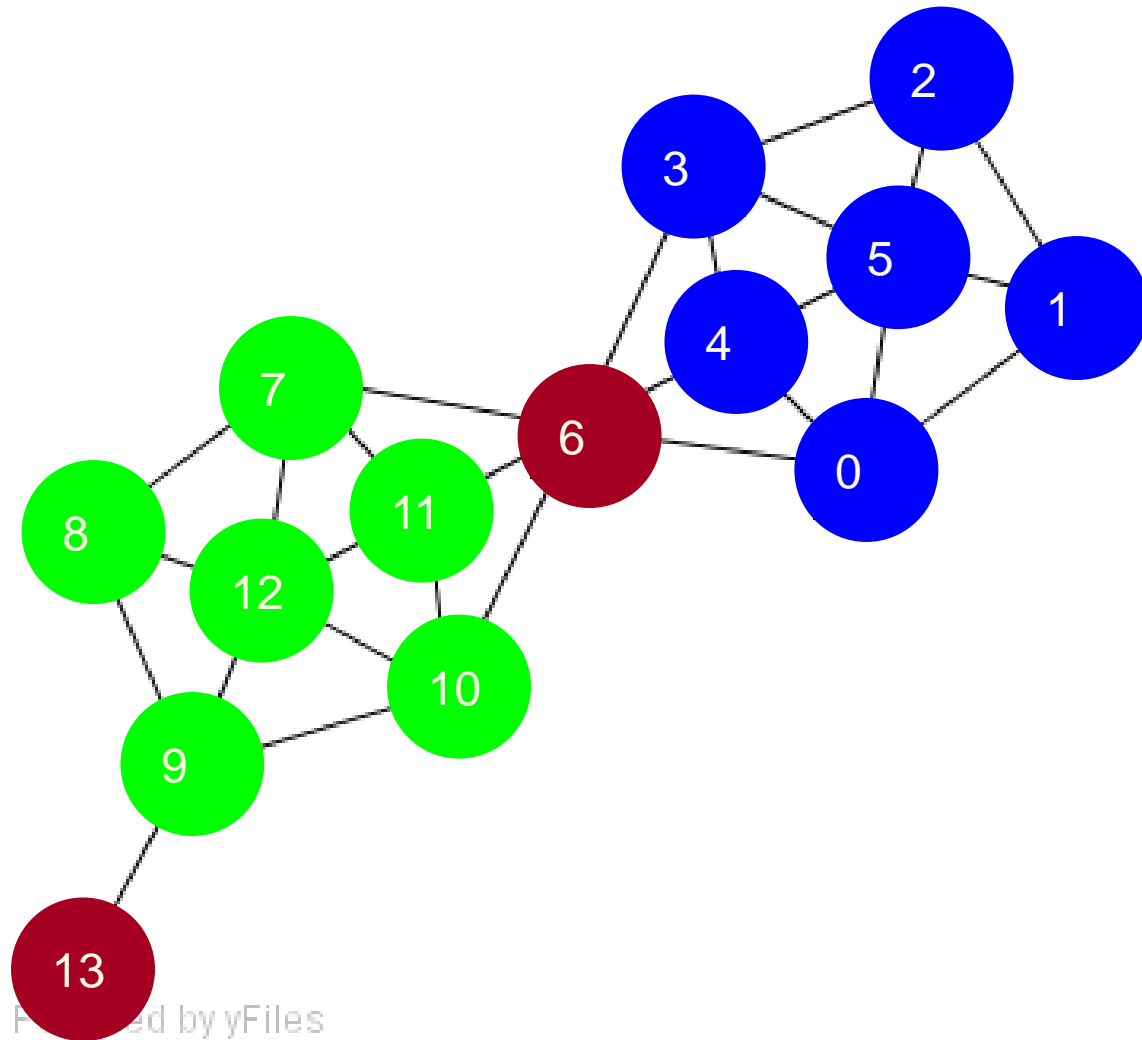
Algorithm

$$\mu = 2$$
$$\varepsilon = 0.7$$



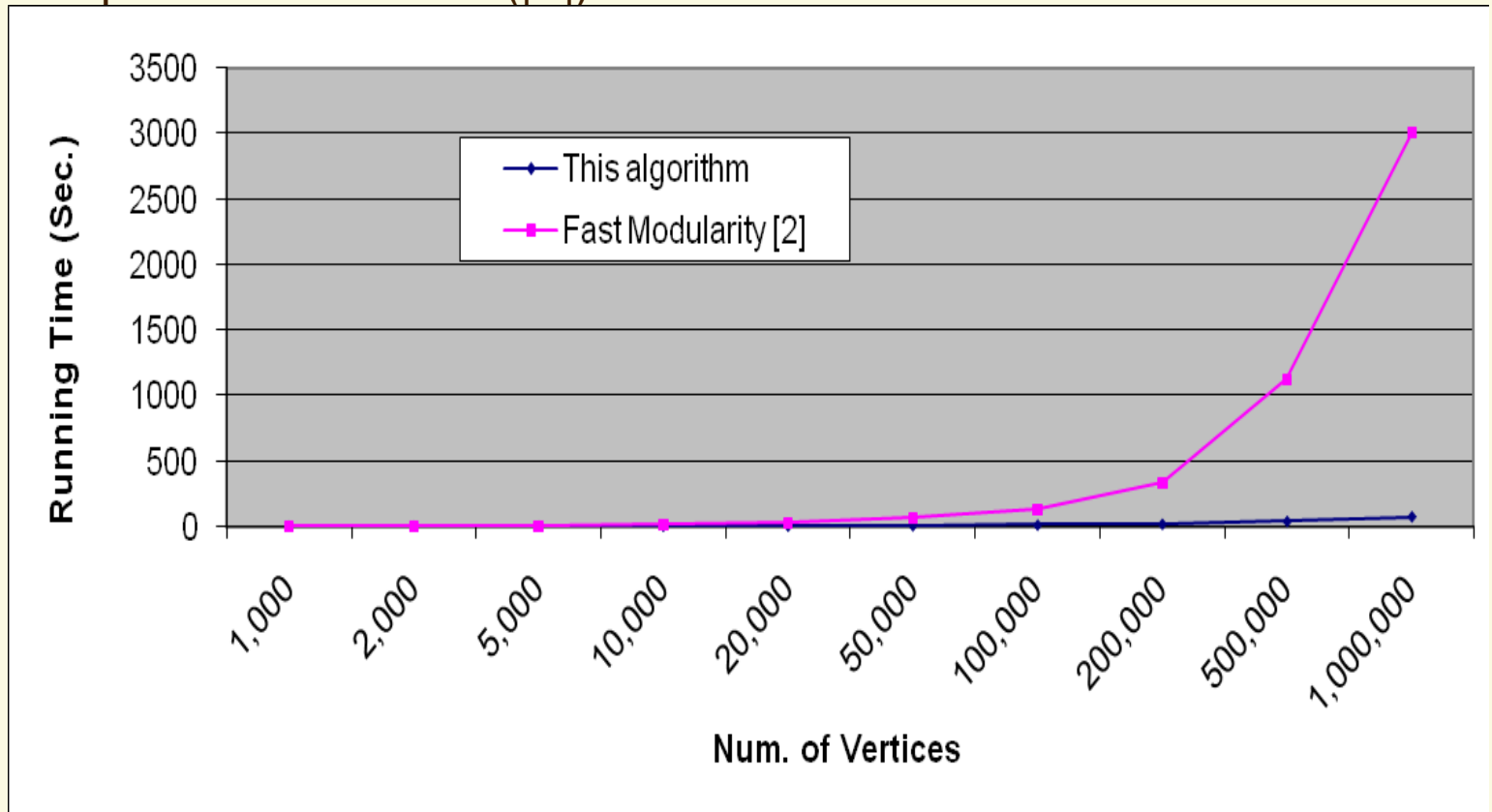
Algorithm

$$\mu = 2$$
$$\varepsilon = 0.7$$



Running Time

- ✓ Running time = $O(|E|)$
- ✓ For sparse networks = $O(|V|)$



[2] A. Clauset, M. E. J. Newman, & C. Moore, *Phys. Rev. E* **70**, 066111 (2004).

Summary

- ✓ Cluster analysis groups objects based on their similarity and has wide applications
- ✓ Measure of similarity can be computed for various types of data
- ✓ Clustering algorithms can be categorized into partitioning methods, hierarchical methods, density-based methods, grid-based methods, and model-based methods
- ✓ K-means and K-medoids algorithms are popular partitioning-based algorithms
- ✓ AGNES and Diana are interesting hierarchical clustering algorithms, and there are also probabilistic hierarchical clustering algorithms
- ✓ DBSCAN, OPTICS, and DENCLU are interesting density-based algorithms
- ✓ Quality of clustering results can be evaluated in various ways
- ✓ Graph Clustering:
 - min-cut vs. sparsest cut
 - High-dimensional clustering methods
 - Graph-specific clustering methods, e.g., SCAN

Papers to read

- ✓ S. Arora, S. Rao, and U. Vazirani. Expander flows, geometric embeddings and graph partitioning. *J. ACM*, 56:5:1–5:37, 2009.
- ✓ I. Davidson and S. S. Ravi. Clustering with constraints: Feasibility issues and the k-means algorithm. *SDM'05*
- ✓ I. Davidson, K. L. Wagstaff, and S. Basu. Measuring constraint-set utility for partitional clustering algorithms. *PKDD'06*
- ✓ C. Fraley and A. E. Raftery. Model-based clustering, discriminant analysis, and density estimation. *J. American Stat. Assoc.*, 97:611–631, 2002.
- ✓ G. Jeh and J. Widom. SimRank: a measure of structural-context similarity. *KDD'02*
- ✓ H.-P. Kriegel, P. Kroeger, and A. Zimek. Clustering high dimensional data: A survey on subspace clustering, pattern-based clustering, and correlation clustering. *ACM Trans. Knowledge Discovery from Data (TKDD)*, 3, 2009.

Papers to read

- ✓ A. Y. Ng, M. I. Jordan, and Y. Weiss. On spectral clustering: Analysis and an algorithm. *NIPS'01*
- ✓ J. Pei, X. Zhang, M. Cho, H. Wang, and P. S. Yu. Maple: A fast algorithm for maximal pattern-based clustering. *ICDM'03*
- ✓ M. Radovanović, A. Nanopoulos, and M. Ivanović. Nearest neighbors in high-dimensional data: the emergence and influence of hubs. *ICML'09*
- ✓ S. E. Schaeffer. Graph clustering. *Computer Science Review*, 1:27–64, 2007.
- ✓ A. K. H. Tung, J. Han, L. V. S. Lakshmanan, and R. T. Ng. Constraint-based clustering in large databases. *ICDT'01*
- ✓ A. Tanay, R. Sharan, and R. Shamir. Biclustering algorithms: A survey. In *Handbook of Computational Molecular Biology*, Chapman & Hall, 2004.
- ✓ H. Wang, W. Wang, J. Yang, and P. S. Yu. Clustering by pattern similarity in large data sets. *SIGMOD'02*
- ✓ X. Xu, N. Yuruk, Z. Feng, and T. A. J. Schweiger. SCAN: A structural clustering algorithm for networks. *KDD'07*

HAPPY THANKSGIVING!

