



# Clustering

Dr. N. Abdolvand

Data Mining

## Data Mining: Concepts and Techniques

(3<sup>rd</sup> ed.)

### — Chapter 10 —

Jiawei Han, Micheline Kamber, and Jian Pei  
University of Illinois at Urbana-Champaign &  
Simon Fraser University

©2013 Han, Kamber & Pei. All rights reserved.

## Chapter 10. Cluster Analysis: Basic Concepts and Methods

- ◆ Cluster Analysis: Basic Concepts
- ◆ Partitioning Methods
- ◆ Hierarchical Methods
- ◆ Density-Based Methods
- ◆ Grid-Based Methods
- ◆ Evaluation of Clustering
- ◆ Summary

3

## 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

4

## Clustering: Application Examples

- ◆ Biology: taxonomy of living things: kingdom, phylum, class, order, family, genus and species
- ◆ Information retrieval: document clustering
- ◆ 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
- ◆ Economic Science: market research

5

## 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

6

## Quality: 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

7

## 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:
  - There is 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

8

## 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, CLARANS
- ◆ Hierarchical approach:
  - Create a hierarchical decomposition of the set of data (or objects) using some criterion
  - 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

9

## 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

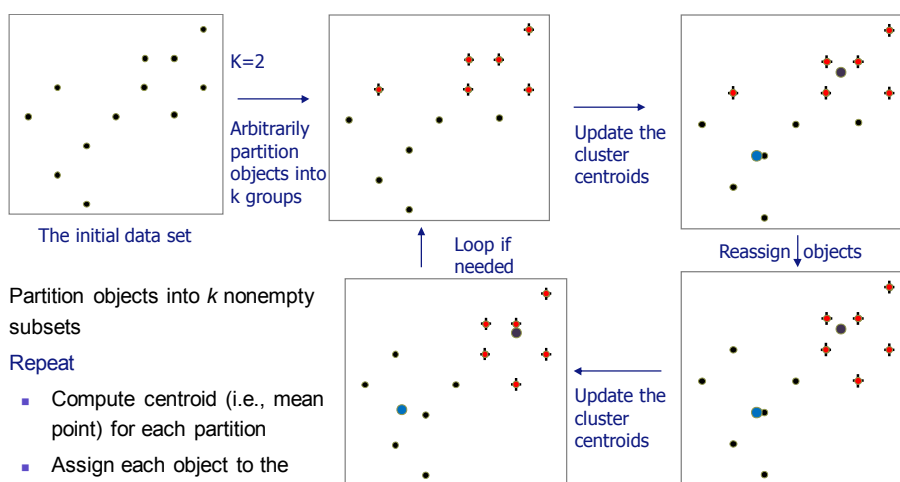
10

## 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

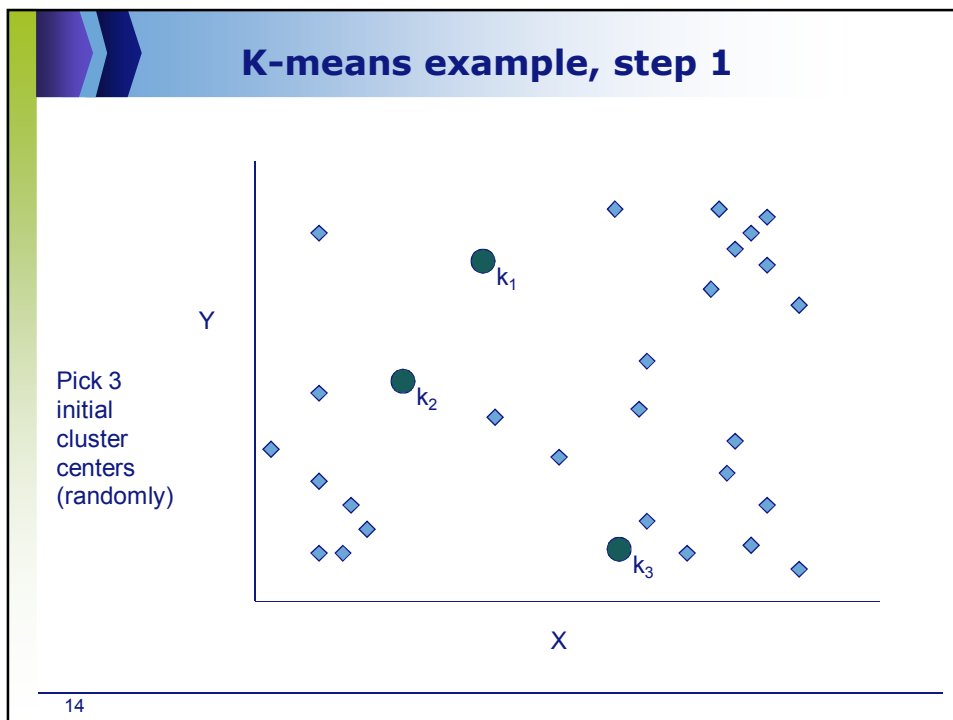
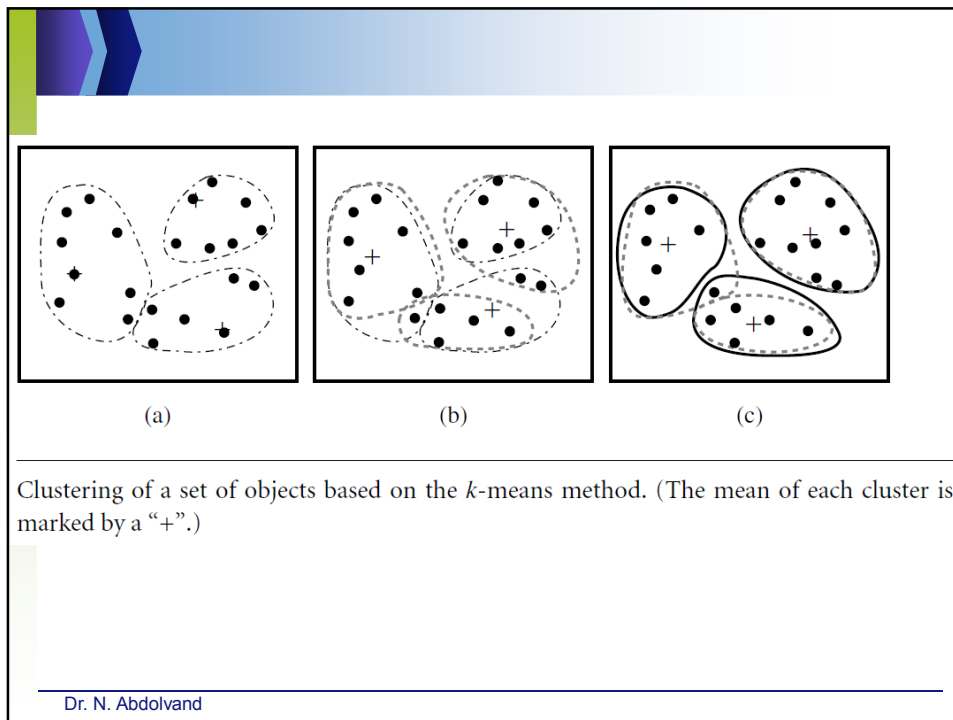
11

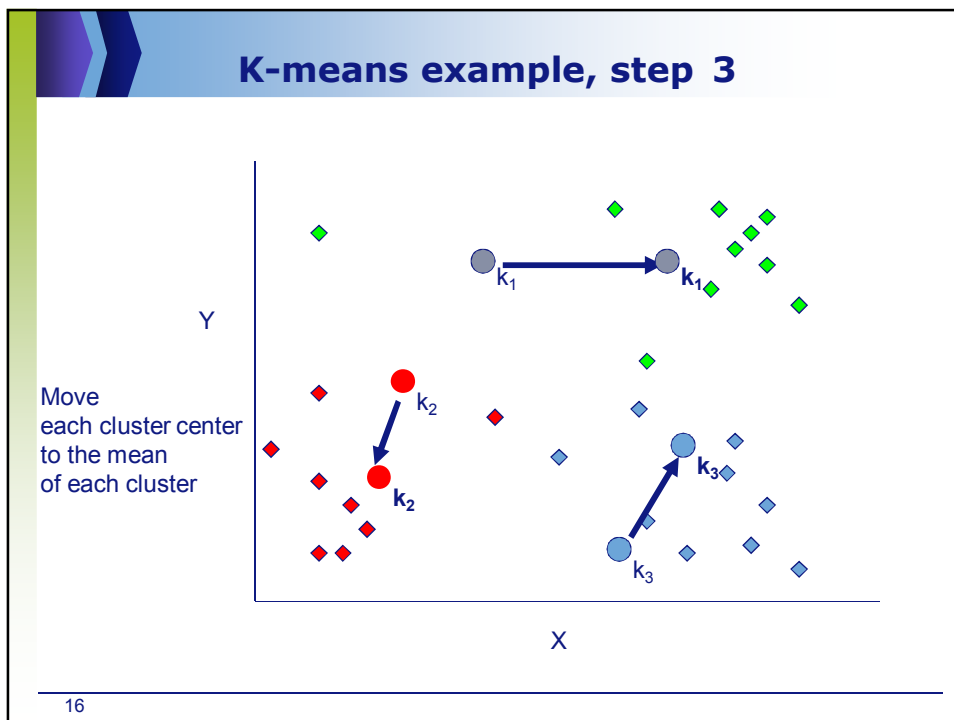
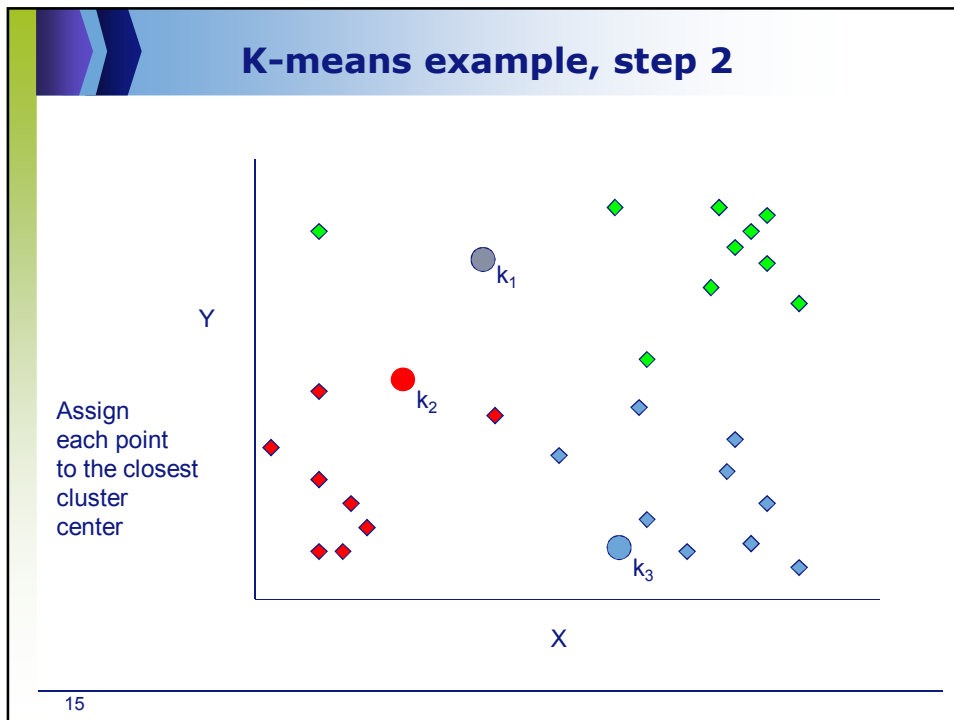
## An Example of *K*-Means Clustering



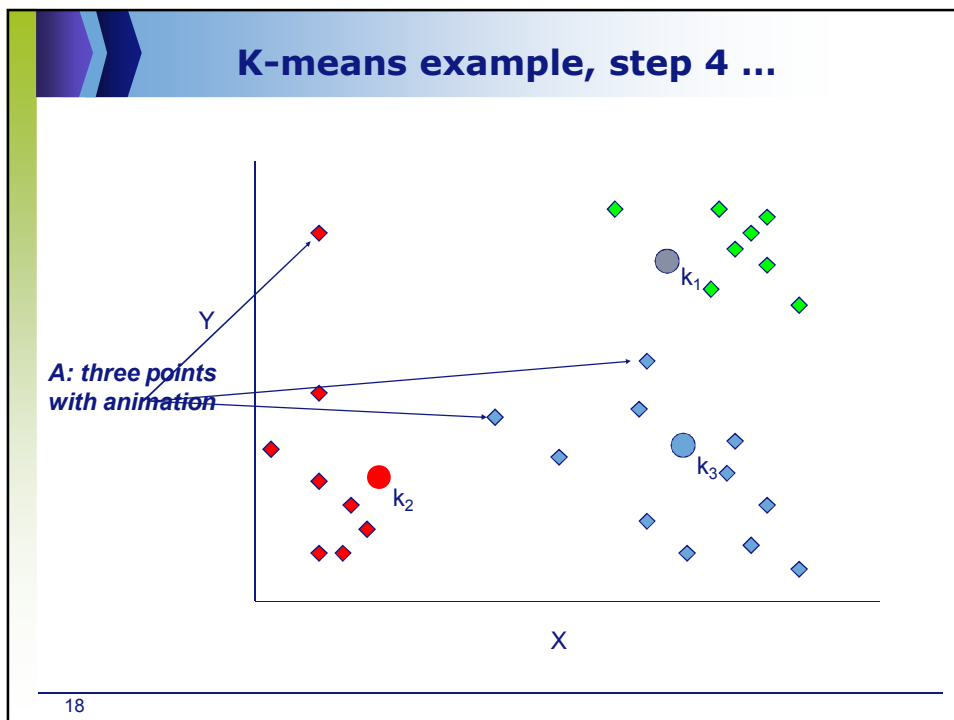
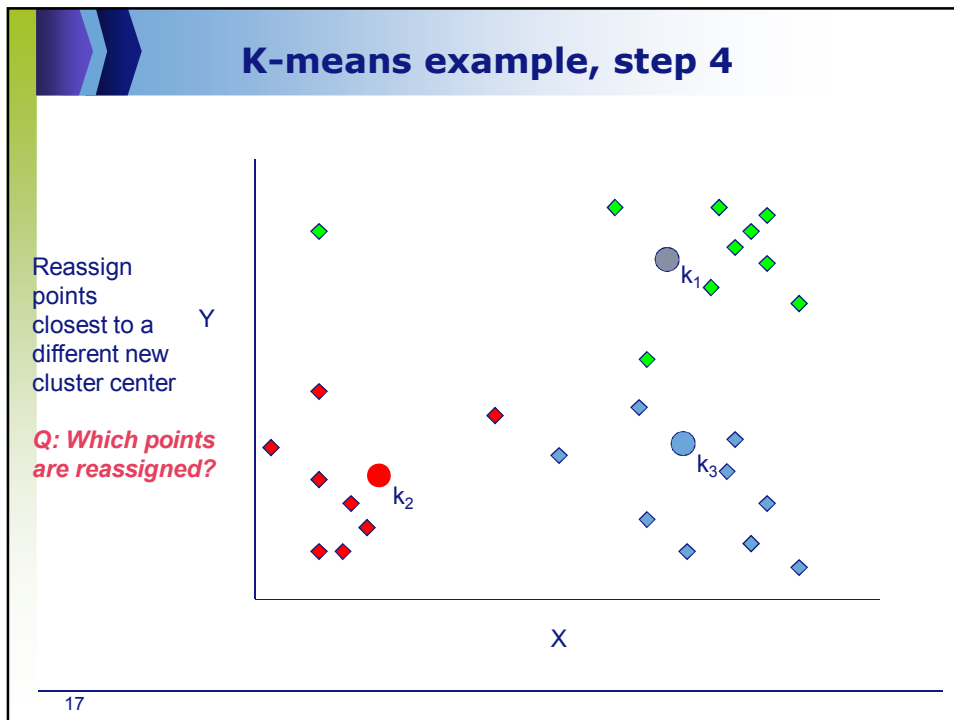
- Partition objects into  $k$  nonempty subsets
- Repeat
  - Compute centroid (i.e., mean point) for each partition
  - Assign each object to the cluster of its nearest centroid
- Until no change

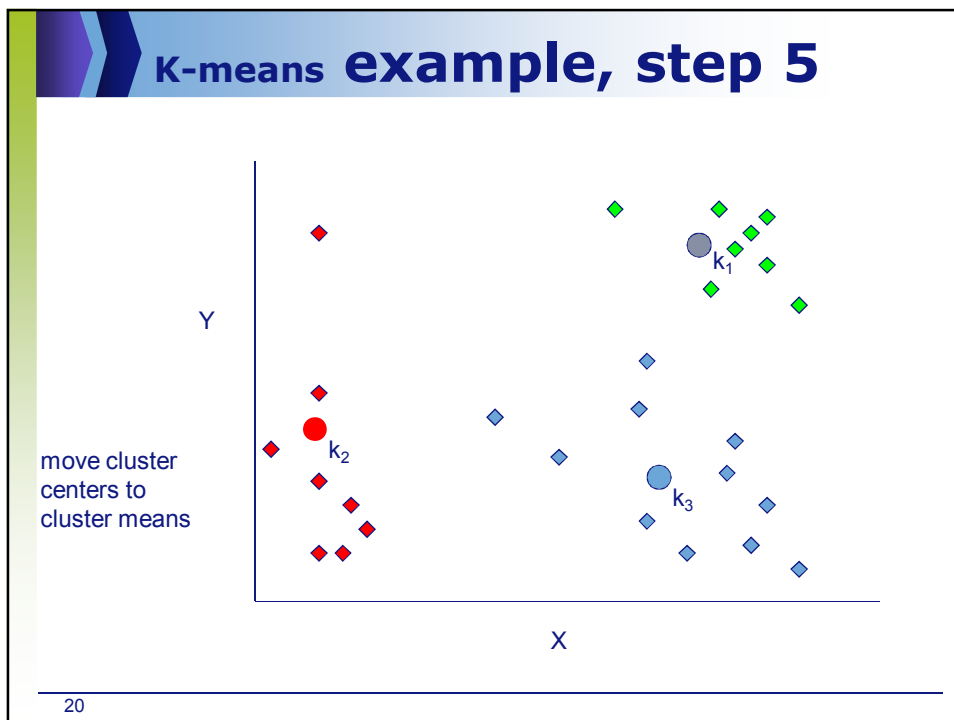
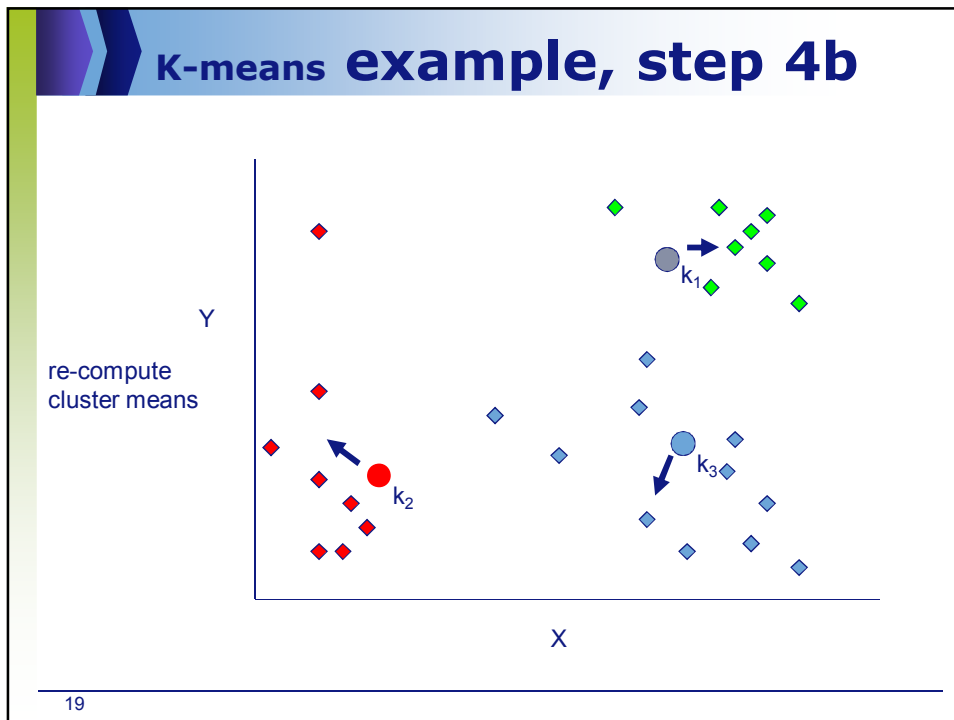
12











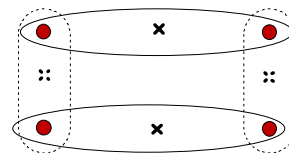
## Comments on the *K-Means* Method

- ◆ **Strength:** *Efficient:*  $O(tkn)$ , where  $n$  is # objects,  $k$  is # clusters, and  $t$  is # iterations. Normally,  $k, t \ll n$ .
  - Comparing: PAM:  $O(k(n-k)^2)$ , CLARA:  $O(ks^2 + k(n-k))$
- ◆ **Comment:** Often terminates at a *local optimal*
- ◆ **Weakness**
  - Applicable only to objects in a continuous  $n$ -dimensional space
    - Using the  $k$ -modes method for categorical data
    - In comparison,  $k$ -medoids can be applied to a wide range of data
  - Need to specify  $k$ , the *number* of clusters, in advance (there are ways to automatically determine the best  $k$  (see Hastie et al., 2009))
  - Sensitive to noisy data and *outliers*
  - Not suitable to discover clusters with *non-convex shapes*

21

## Variations of the *K-Means* Method

- ◆ Most of the variants of the *k-means* which differ in
  - Selection of the initial  $k$  means
  - Dissimilarity calculations
  - Strategies to calculate cluster means
- ◆ Handling categorical data: *k-modes*
  - Replacing means of clusters with modes
  - Using new dissimilarity measures to deal with categorical objects
  - Using a frequency-based method to update modes of clusters
  - A mixture of categorical and numerical data: *k-prototype* method



22

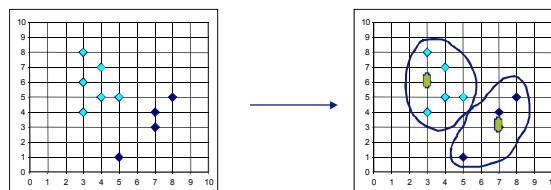
## K-means variations

- ◆ **K-medoids** – instead of mean, use medians of each cluster
  - Mean of 1, 3, 5, 7, 9 is **5**
  - Mean of 1, 3, 5, 7, 1009 is **205**
  - Median of 1, 3, 5, 7, 1009 is **5**
  - Median advantage: not affected by extreme values
- ◆ For large databases, use sampling

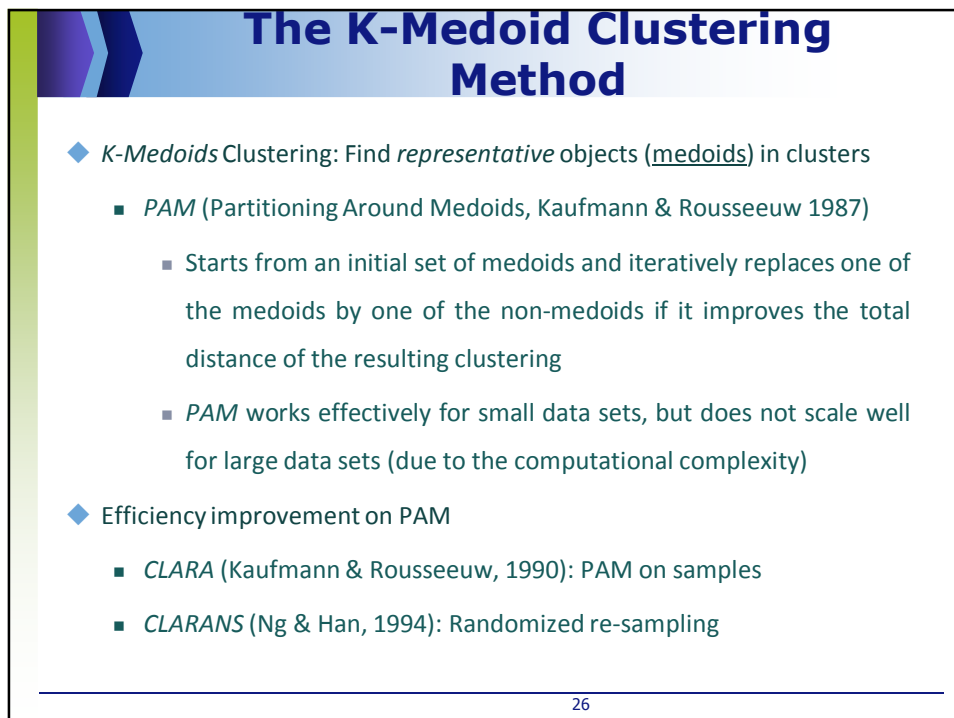
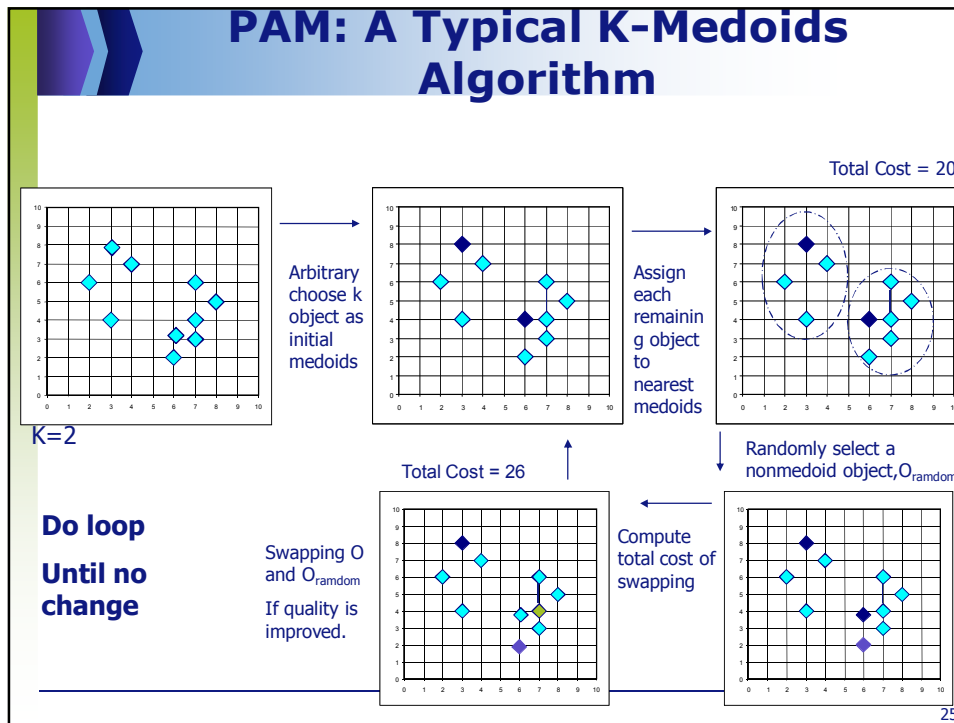
23

## 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

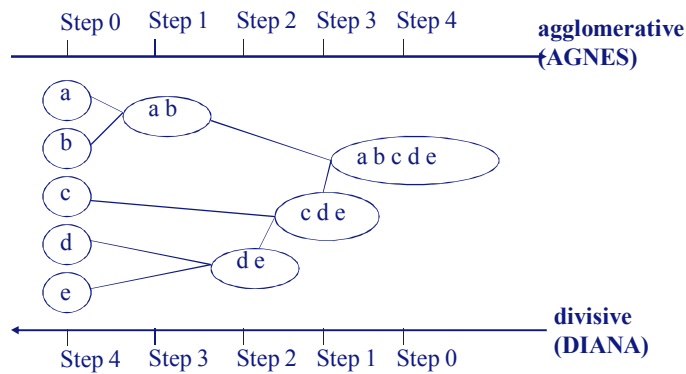


24



## Hierarchical Clustering

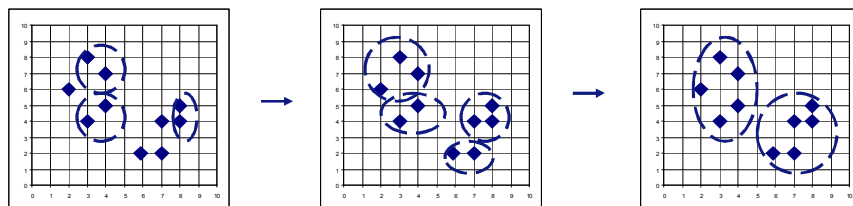
- ◆ Use distance matrix as clustering criteria. This method does not require the number of clusters  $k$  as an input, but needs a termination condition



27

## 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

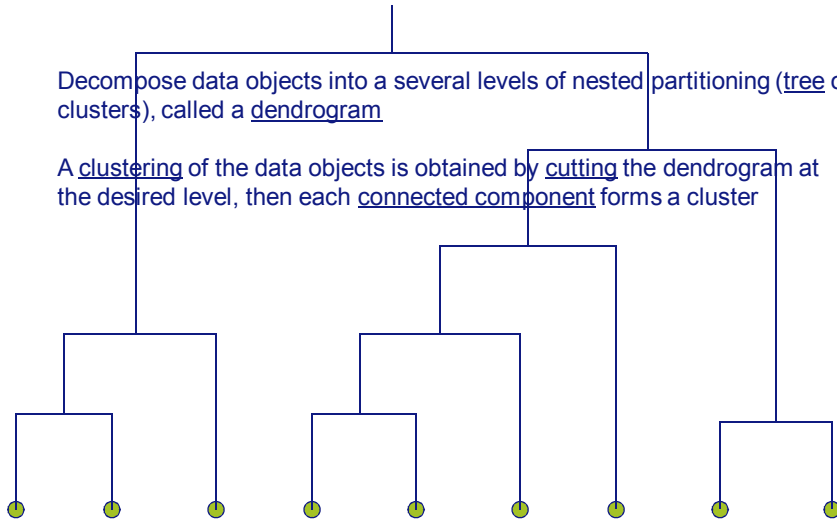


28

## Dendrogram: Shows How Clusters are Merged

Decompose data objects into a several levels of nested partitioning (tree of clusters), called a dendrogram

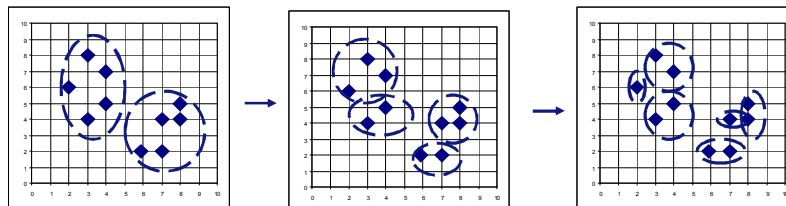
A clustering of the data objects is obtained by cutting the dendrogram at the desired level, then each connected component forms a cluster



29

## DIANA (Divisive Analysis)

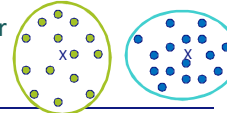
- ◆ Introduced in Kaufmann and Rousseeuw (1990)
- ◆ Implemented in statistical analysis packages, e.g., Splus
- ◆ Inverse order of AGNES
- ◆ Eventually each node forms a cluster on its own



30

## Distance between Clusters

- ◆ Single link: smallest distance between an element in one cluster and an element in the other, i.e.,  $\text{dist}(K_i, K_j) = \min(t_{ip}, t_{jq})$
- ◆ Complete link: largest distance between an element in one cluster and an element in the other, i.e.,  $\text{dist}(K_i, K_j) = \max(t_{ip}, t_{jq})$
- ◆ Average: avg distance between an element in one cluster and an element in the other, i.e.,  $\text{dist}(K_i, K_j) = \text{avg}(t_{ip}, t_{jq})$
- ◆ Centroid: distance between the centroids of two clusters, i.e.,  $\text{dist}(K_i, K_j) = \text{dist}(C_i, C_j)$
- ◆ Medoid: distance between the medoids of two clusters, i.e.,  $\text{dist}(K_i, K_j) = \text{dist}(M_i, M_j)$ 
  - Medoid: a chosen, centrally located object in the cluster



31

## Centroid, Radius and Diameter of a Cluster (for numerical data sets)

- ◆ Centroid: the “middle” of a cluster

$$C_m = \frac{\sum_{i=1}^N (t_{ip})}{N}$$

- ◆ Radius: square root of average distance from any point of the cluster to its centroid

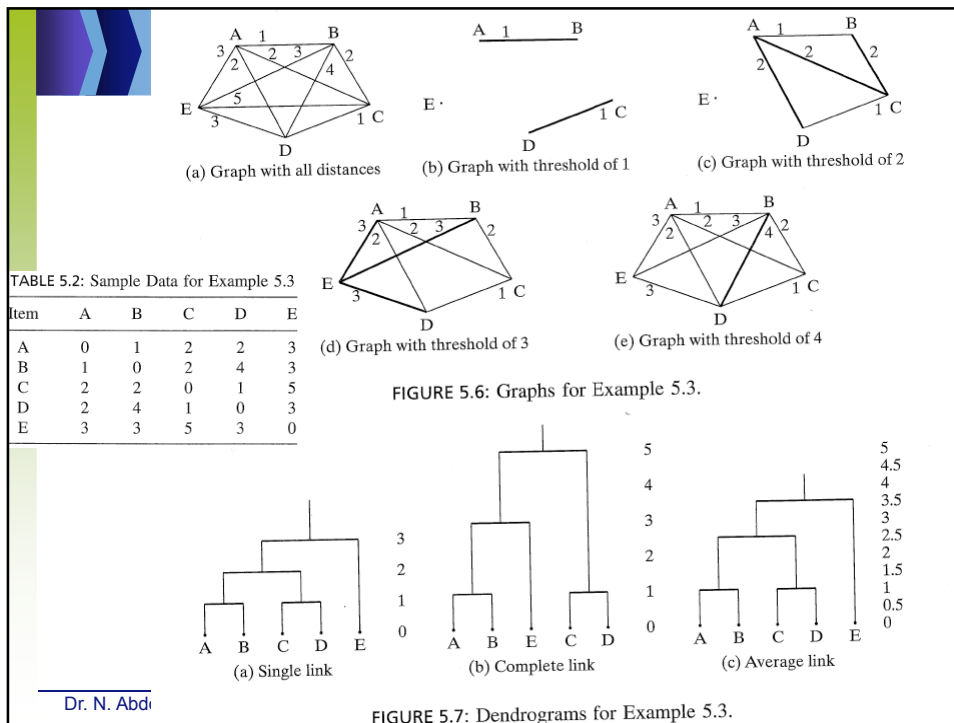
$$R_m = \sqrt{\frac{\sum_{i=1}^N (t_{ip} - c_m)^2}{N}}$$

- ◆ Diameter: square root of average mean squared distance between all pairs of points in the cluster

$$D_m = \sqrt{\frac{\sum_{i=1}^N \sum_{j=1}^N (t_{ip} - t_{jq})^2}{N(N-1)}}$$

32





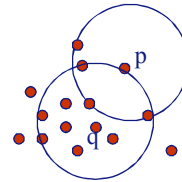
## Density-Based Clustering Methods

- ◆ 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 neighbourhood
  - *MinPts*: Minimum number of points in an *Eps*-neighbourhood 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  $N_{Eps}(q)$
  - core point condition:

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



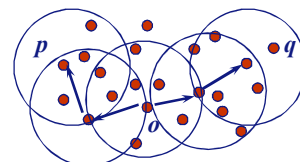
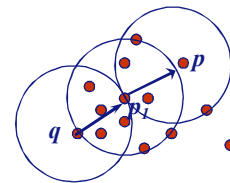
MinPts = 5

Eps = 1 cm

35

## Density-Reachable and Density-Connected

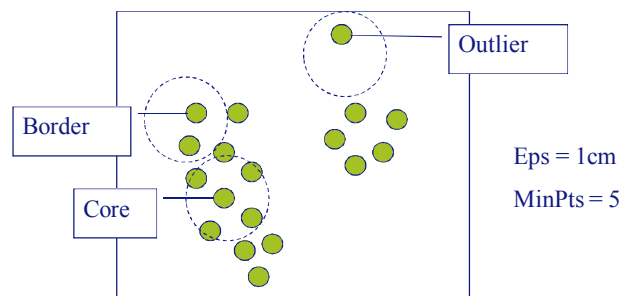
- ◆ Density-reachable:
  - A point  $p$  is **density-reachable** from a point  $q$  w.r.t. *Eps*, *MinPts* if there is a chain of points  $p_1, \dots, p_n$ ,  $p_1 = q$ ,  $p_n = p$  such that  $p_{i+1}$  is directly density-reachable from  $p_i$
- ◆ Density-connected
  - A point  $p$  is **density-connected**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



36

## 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



37

## DBSCAN: The Algorithm

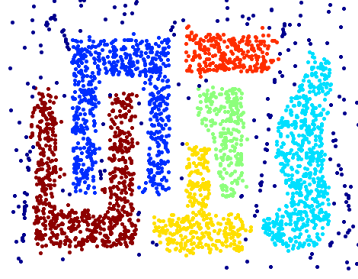
- ◆ Arbitrary select a point  $p$
- ◆ Retrieve all points density-reachable from  $p$  w.r.t.  $Eps$  and  $MinPts$
- ◆ If  $p$  is a core point, a cluster is formed
- ◆ If  $p$  is a border point, no points are density-reachable from  $p$  and DBSCAN visits the next point of the database
- ◆ Continue the process until all of the points have been processed
- ◆ If a spatial index is used, the computational complexity of DBSCAN is  $O(n \log n)$ , where  $n$  is the number of database objects. Otherwise, the complexity is  $O(n^2)$

38

## When DBSCAN Works Well



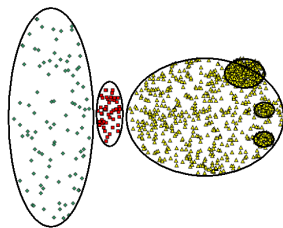
Original Points



Clusters

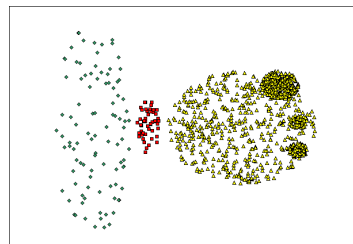
- Resistant to Noise
- Can handle clusters of different shapes and sizes

## When DBSCAN Does NOT Work Well

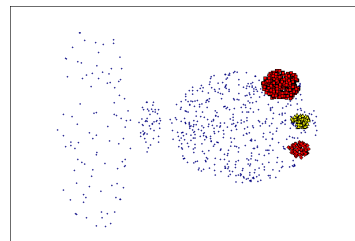


Original Points

- Cannot handle Varying densities
- sensitive to parameters



(MinPts=4, Eps=9.92).



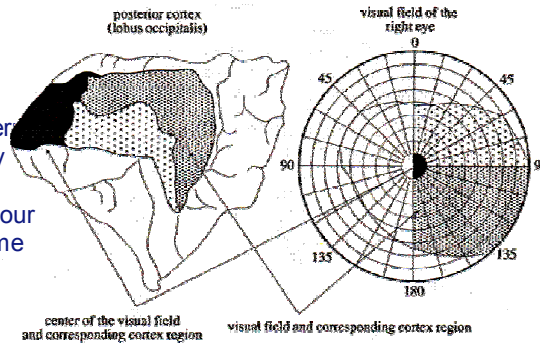
(MinPts=4, Eps=9.75)

## Brain's self-organization

The brain maps the external multidimensional representation of the world into a similar 1 or 2 - dimensional internal representation.

That is, the brain processes the external signals in a topology-preserving way

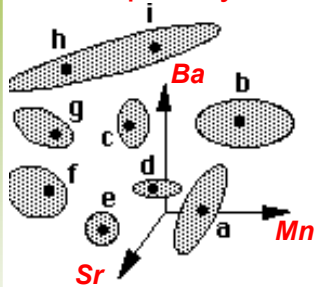
Mimicking the way the brain learns, our system should be able to do the same thing.



f1

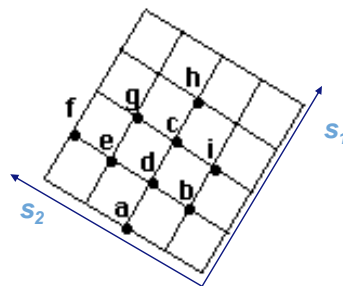
## Concept of the SOM.

Input space  
Input layer



Cluster centers (code vectors)

Reduced feature space  
Map layer



Place of these code vectors in the reduced space

**Clustering and ordering** of the cluster centers in a two dimensional grid

f2

## Self-Organizing Feature Map (SOM)

- ◆ SOMs, also called topological ordered maps, or Kohonen Self-Organizing Feature Map (KSOMs)
- ◆ It maps all the points in a high-dimensional source space into a 2 to 3-d target space, s.t., the distance and proximity relationship (i.e., topology) are preserved as much as possible
- ◆ Similar to k-means: cluster centers tend to lie in a low-dimensional manifold in the feature space
- ◆ Clustering is performed by having several units competing for the current object
  - The unit whose weight vector is closest to the current object wins
  - The winner and its neighbors learn by having their weights adjusted
- ◆ SOMs are believed to resemble processing that can occur in the brain
- ◆ Useful for visualizing high-dimensional data in 2- or 3-D space

May 15, 2017

43

Data Mining: Concepts and  
Techniques

## Self-Organizing Feature Map (SOM)

- ◆ SOMs, also called topological ordered maps, or Kohonen Self-Organizing Feature Map (KSOMs)
- ◆ It maps all the points in a high-dimensional source space into a 2 to 3-d target space, s.t., the distance and proximity relationship (i.e., topology) are preserved as much as possible
- ◆ Similar to k-means: cluster centers tend to lie in a low-dimensional manifold in the feature space
- ◆ Clustering is performed by having several units competing for the current object
  - The unit whose weight vector is closest to the current object wins
  - The winner and its neighbors learn by having their weights adjusted
- ◆ SOMs are believed to resemble processing that can occur in the brain
- ◆ Useful for visualizing high-dimensional data in 2- or 3-D space

May 15, 2017

44

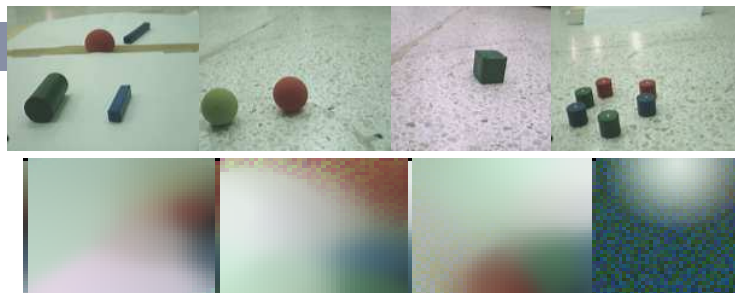
Data Mining: Concepts and  
Techniques

### Examples-1

network. We use a fixed learning rate of  $\eta=1.0E-4$  and 1000 epochs. About 200 images were used in order to extract the pixel values for training.

- Some of the original images and unsuccessful & successful colour maps are shown below:

#### Examples



### Examples-2

- **Semantic Maps:** A useful method of visualisation of the SOM structure achieved at the end of training assigns class labels in a 2D lattice depending on how each test pattern (not seen before) excites a particular neuron.

#### Examples

- The neurons in the lattice are partitioned to a number of *coherent regions*, coherent in the sense that each grouping of neurons represents a distinct set of contiguous symbols or labels.
- An example is shown below, where we assume that we have trained the map for 16 different animals.
- We use a lattice of 10x10 output neurons.

## Examples-3

dog	dog	fox	fox	fox	cat	cat	cat	eagle	eagle
dog	dog	fox	fox	fox	cat	cat	cat	eagle	eagle
wolf	wolf	wolf	fox	cat	tiger	tiger	tiger	owl	owl
wolf	wolf	lion	lion	lion	tiger	tiger	tiger	hawk	hawk
wolf	wolf	lion	lion	lion	tiger	tiger	tiger	hawk	hawk
wolf	wolf	lion	lion	lion	owl	dove	hawk	dove	dove
horse	horse	lion	lion	lion	dove	hen	hen	dove	dove
horse	horse	zebra	cow	cow	cow	hen	hen	dove	dove
zebra	zebra	zebra	cow	cow	cow	hen	hen	duck	goose
zebra	zebra	zebra	cow	cow	cow	duck	duck	duck	goose

## Examples

- We observe that there are three distinct clusters of animals: "birds", "peaceful species" and "hunters".

## Measuring Clustering Quality

- ◆ 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



## Evaluation

### ◆ Sum-of-Squared-Error Criterion

- Summing over the squared distances between the clustering objects and their cluster representatives (i.e. the respective cluster centroids) is a standard cost

$$E(C) = \sum_{i=1}^k \sum_{o \in C_i} d(o, cen_i)^2$$

### ◆ Silhouette Value

- Good clusters are those where the genes are close to each other compared to their next closest cluster.

$$b(i) = \min(\text{AVGD\_BETWEEN}(i, k))$$

$$a(i) = \text{AVGD\_WITHIN}(i)$$

$$s(i) = \frac{b(i) - a(i)}{\max(a(i), b(i))}$$

Dr. N. Abdolvand

## 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 clustering algorithms
- ◆ Birch and Chameleon are interesting hierarchical clustering algorithms, and there are also probabilistic hierarchical clustering algorithms
- ◆ DBSCAN, OPTICS, and DENCLU are interesting density-based algorithms
- ◆ STING and CLIQUE are grid-based methods, where CLIQUE is also a subspace clustering algorithm
- ◆ Quality of clustering results can be evaluated in various ways

## References (1)

- ◆ R. Agrawal, J. Gehrke, D. Gunopulos, and P. Raghavan. Automatic subspace clustering of high dimensional data for data mining applications. SIGMOD'98
- ◆ M. R. Anderberg. Cluster Analysis for Applications. Academic Press, 1973.
- ◆ M. Ankerst, M. Breunig, H.-P. Kriegel, and J. Sander. Optics: Ordering points to identify the clustering structure, SIGMOD'99.
- ◆ Beil F., Ester M., Xu X.: "Frequent Term-Based Text Clustering", KDD'02
- ◆ M. M. Breunig, H.-P. Kriegel, R. Ng, J. Sander. LOF: Identifying Density-Based Local Outliers. SIGMOD 2000.
- ◆ M. Ester, H.-P. Kriegel, J. Sander, and X. Xu. A density-based algorithm for discovering clusters in large spatial databases. KDD'96.
- ◆ M. Ester, H.-P. Kriegel, and X. Xu. Knowledge discovery in large spatial databases: Focusing techniques for efficient class identification. SSD'95.
- ◆ D. Fisher. Knowledge acquisition via incremental conceptual clustering. Machine Learning, 2:139-172, 1987.
- ◆ D. Gibson, J. Kleinberg, and P. Raghavan. Clustering categorical data: An approach based on dynamic systems. VLDB'98.
- ◆ V. Ganti, J. Gehrke, R. Ramakrishan. CACTUS Clustering Categorical Data Using Summaries. KDD'99.

52

## References (2)

- ◆ D. Gibson, J. Kleinberg, and P. Raghavan. Clustering categorical data: An approach based on dynamic systems. In Proc. VLDB'98.
- ◆ S. Guha, R. Rastogi, and K. Shim. Cure: An efficient clustering algorithm for large databases. SIGMOD'98.
- ◆ S. Guha, R. Rastogi, and K. Shim. ROCK: A robust clustering algorithm for categorical attributes. In *ICDE'99*, pp. 512-521, Sydney, Australia, March 1999.
- ◆ A. Hinneburg, D.I A. Keim: An Efficient Approach to Clustering in Large Multimedia Databases with Noise. KDD'98.
- ◆ A. K. Jain and R. C. Dubes. Algorithms for Clustering Data. Printice Hall, 1988.
- ◆ G. Karypis, E.-H. Han, and V. Kumar. CHAMELEON: A Hierarchical Clustering Algorithm Using Dynamic Modeling. *COMPUTER*, 32(8): 68-75, 1999.
- ◆ L. Kaufman and P. J. Rousseeuw. Finding Groups in Data: an Introduction to Cluster Analysis. John Wiley & Sons, 1990.
- ◆ E. Knorr and R. Ng. Algorithms for mining distance-based outliers in large datasets. VLDB'98.

53

## References (3)

- ◆ G. J. McLachlan and K.E. Bksford. Mixture Models: Inference and Applications to Clustering. John Wiley and Sons, 1988.
- ◆ R. Ng and J. Han. Efficient and effective clustering method for spatial data mining. VLDB'94.
- ◆ L. Parsons, E. Haque and H. Liu, Subspace Clustering for High Dimensional Data: A Review, SIGKDD Explorations, 6(1), June 2004
- ◆ E. Schikuta. Grid clustering: An efficient hierarchical clustering method for very large data sets. Proc. 1996 Int. Conf. on Pattern Recognition,.
- ◆ G. Sheikholeslami, S. Chatterjee, and A. Zhang. WaveCluster: A multi-resolution clustering approach for very large spatial databases. VLDB'98.
- ◆ A. K. H. Tung, J. Han, L. V. S. Lakshmanan, and R. T. Ng. Constraint-Based Clustering in Large Databases, ICDT'01.
- ◆ A. K. H. Tung, J. Hou, and J. Han. Spatial Clustering in the Presence of Obstacles, ICDE'01
- ◆ H. Wang, W. Wang, J. Yang, and P.S. Yu. Clustering by pattern similarity in large data sets, SIGMOD' 02.
- ◆ W. Wang, Yang, R. Muntz, STING: A Statistical Information grid Approach to Spatial Data Mining, VLDB'97.
- ◆ T. Zhang, R. Ramakrishnan, and M. Livny. BIRCH : An efficient data clustering method for very large databases. SIGMOD'96.
- ◆ Xiaoxin Yin, Jiawei Han, and Philip Yu, "[LinkClus: Efficient Clustering via Heterogeneous Semantic Links](#)", in Proc. 2006 Int. Conf. on Very Large Data Bases (VLDB'06), Seoul, Korea, Sept. 2006.



# Questions?