

clustering

→ an supervised



# General Applications of Clustering

---

- Pattern Recognition
- Spatial Data Analysis (공간데이터)
  - create thematic maps in GIS by clustering feature spaces
  - detect spatial clusters and explain them in spatial data mining
- Image Processing
- Economic Science (especially market research) 문서성향 고객군집화
- WWW 웹문서
  - Document classification
  - Cluster Weblog data to discover groups of similar access patterns



# Examples of Clustering Applications

---

- 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 보험청구액 → 보험사기 감시
- 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 진앙지 표시 → 군집 (지진 분포)



# What Is Good Clustering?

---

- A good clustering method will produce high quality clusters with
  - high intra-class similarity 높은 군집내 유사성
  - low inter-class similarity 낮은 군집간 유사성



# Similarity and Dissimilarity Between Objects

---

- Distances are normally used to measure the similarity or dissimilarity between two data objects

- Some popular ones include: *Minkowski distance*: *minkowski sum*

$$d(i, j) = \sqrt[q]{(|x_{i1} - x_{j1}|^q + |x_{i2} - x_{j2}|^q + \dots + |x_{ip} - x_{jp}|^q)} \quad q=2$$

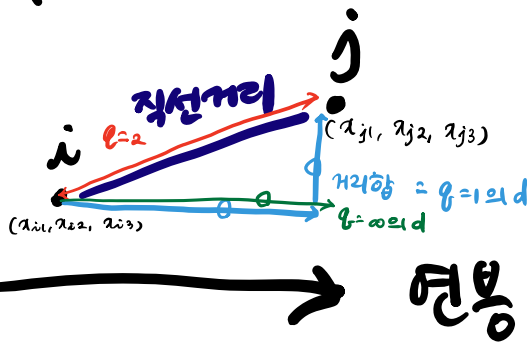
where  $i = (x_{i1}, x_{i2}, \dots, x_{ip})$  and  $j = (x_{j1}, x_{j2}, \dots, x_{jp})$  are two  $p$ -dimensional data objects, and  $q$  is a positive integer

- If  $q = 1$ ,  $d$  is Manhattan distance

$$d(i, j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{ip} - x_{jp}|$$

# Minkowski distance

나이



직선거리만큼이라  
유사도 측정 [0, 1]

아주  
다르다

동일

$q=2$

유클리드  $d = \sqrt{(x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + (x_{i3} - x_{j3})^2}$

$q=1$

$$d = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + |x_{i3} - x_{j3}|$$

$q=\infty$

$$d = \max(|x_{i1} - x_{j1}|, |x_{i2} - x_{j2}|, |x_{i3} - x_{j3}|)$$

→  $\phi$  제일 큰 값이 다가옴.

$\boxed{3^\infty}$  vs  $\cancel{\infty}$

# Similarity and Dissimilarity Between Objects (Cont.)

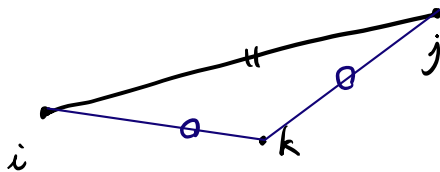
- If  $q = 2$ ,  $d$  is Euclidean distance:

$$d(i, j) = \sqrt{(|x_{i1} - x_{j1}|^2 + |x_{i2} - x_{j2}|^2 + \dots + |x_{ip} - x_{jp}|^2)}$$

Properties

- $d(i, j) \geq 0$
- $d(i, i) = 0$  같은 점거리
- $d(i, j) = d(j, i)$  거리는 방향성 X
- $d(i, j) \leq d(i, k) + d(k, j)$

4가지 특징



+triangular inequality



# Major Clustering Approaches

클러스터링 방법

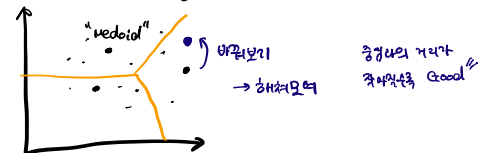
- Partitioning algorithms: Construct various partitions and then evaluate them by some criterion    공간을 쪼개면서
- Hierarchy algorithms: Create a hierarchical decomposition of the set of data (or objects) using some criterion    계층적
- Density-based: based on connectivity and density functions    밀도기반
- Grid-based: based on a multiple-level granularity structure    격자 기반 (밀도 높낮이)
- Model-based: A model is hypothesized for each of the clusters and the idea is to find the best fit of that model to each other    데이터 분포에 적합한 모델

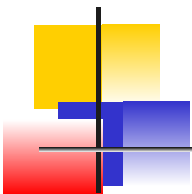


# Partitioning Algorithms: Basic Concept

단점)  $k$ 가 주어지지 않는다는 것

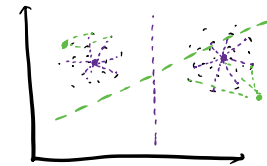
- Partitioning method: Construct a partition of a database  $D$  of  $n$  objects into a set of  $k$  clusters  
(주어진 값)  $n$  data item  
 $k$ 개의 군집
- Given a  $k$ , find a partition of  $k$  clusters that optimizes the chosen partitioning criterion  
나누는 최적  $k$ 를 찾기
- Global optimal: exhaustively enumerate all partitions  
어려움
- Heuristic methods:  $k$ -means and  $k$ -medoids algorithms
- (Spss에 있음)  $k$ -means (MacQueen'67): Each cluster is represented by the center of the cluster  
더 좋음  
 임의의 중심점 선택
- $k$ -medoids or PAM (Partition around medoids) (Kaufman & Rousseeuw'87): Each cluster is represented by one of the objects in the cluster  
실제로 있는 임의의 대표점 선택  
 → 대표할 데이터 지리에 부합함





# The *K-Means* Clustering Method

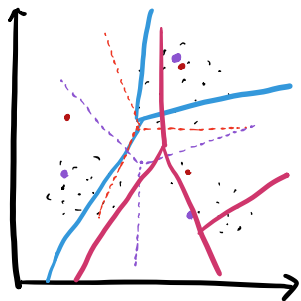
$k$  = 작  $\rightarrow$  큰  
 군집의 품질을 고려해서  
 $k$  값 정하기



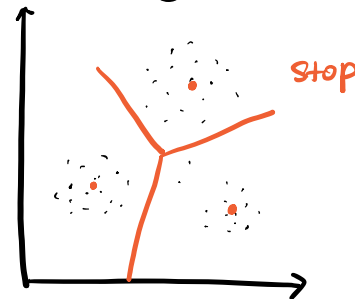
중심점을 거리 계산을 계산.  
 중심점 (seed)까지

주어진  $k$

- Given  $k$ , the  $k$ -means algorithm is implemented in 4 steps:
- Partition objects into  $k$  nonempty subsets
- Compute seed points as the centroids of the clusters of the current partition. The centroid is the center (mean point) of the cluster.
- Assign each object to the cluster with the nearest seed point.
- Go back to Step 2, stop when no more new assignment.



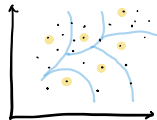
$k=3$  ① 임의의 3개의 공간으로 나눔 —  
 ② 중심좌표 (seed) 찾기 •  
 ③ 중심좌표와 가까운 쪽으로 assign  
 ④ 2-3 반복  $\rightarrow$  change가 없으면 stop —  
 IF  $k=4$ , 이상한 모양



# The *K-Medoids* Clustering Method

- Find *representative* objects, called medoids, in clusters
- *PAM* (Partitioning Around Medoids, 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
- *CLARA* (Kaufmann & Rousseeuw, 1990)

→ sample 축소. ⇒ sample만 가지고 PAM 알고리즘 돌리기



파티셔닝

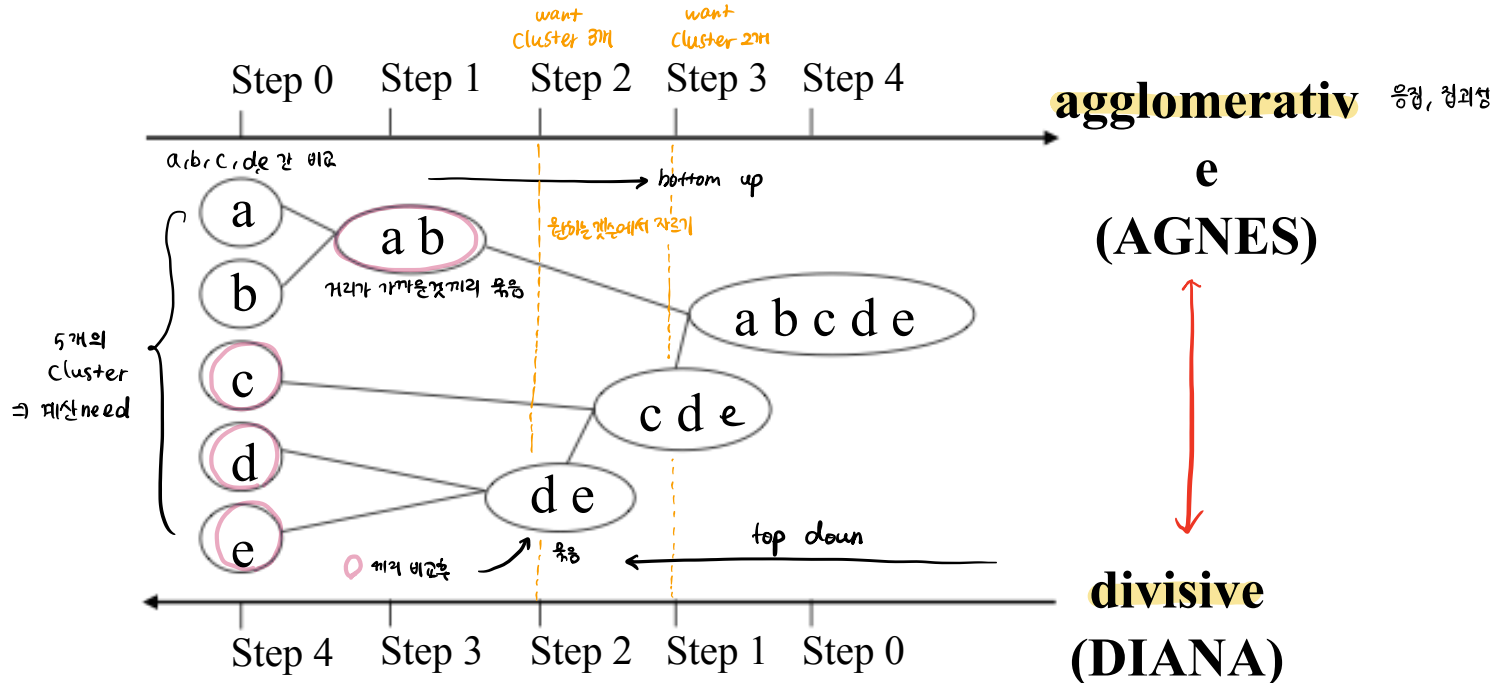
k값 필요 (군집의 수)

임의의 k를 선택 → 데이터 할당 → 중심점 찾고 → 데이터 할당

# Hierarchical Clustering

계층화 (장점:  $k$ 값 필요X)      단점: 데이터 양 → 효율성 ↓

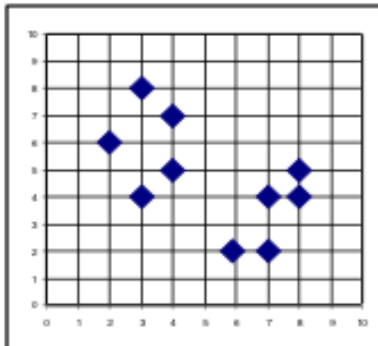
Use distance matrix as clustering criteria. This method does not require the number of clusters  $k$  as an input, but needs a termination condition



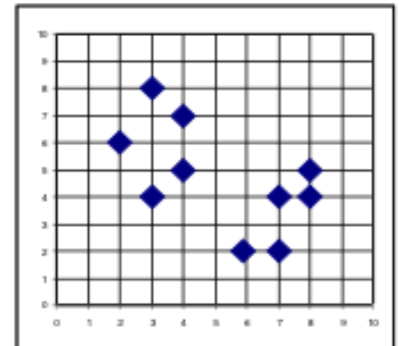
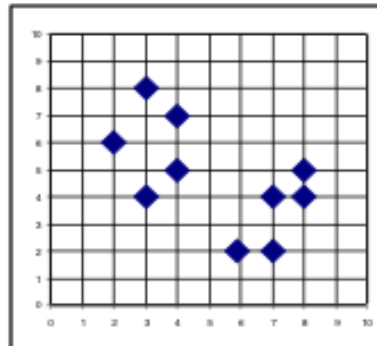
# AGNES (Agglomerative Nesting)

- Introduced in Kaufmann and Rousseeuw (1990)
- Implemented in statistical analysis 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

점점



점점 need

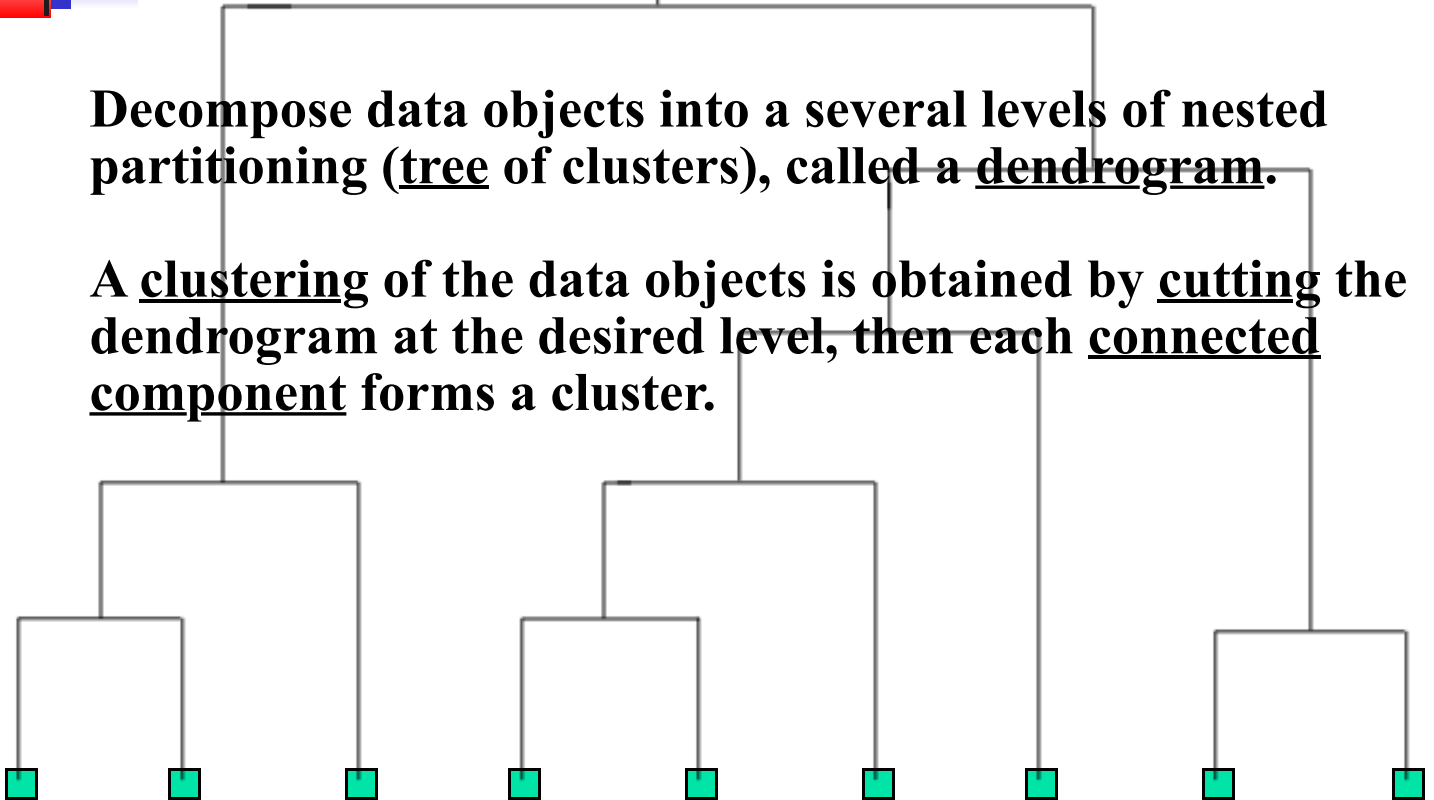




# ***A Dendrogram Shows How the Clusters are Merged Hierarchically***

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

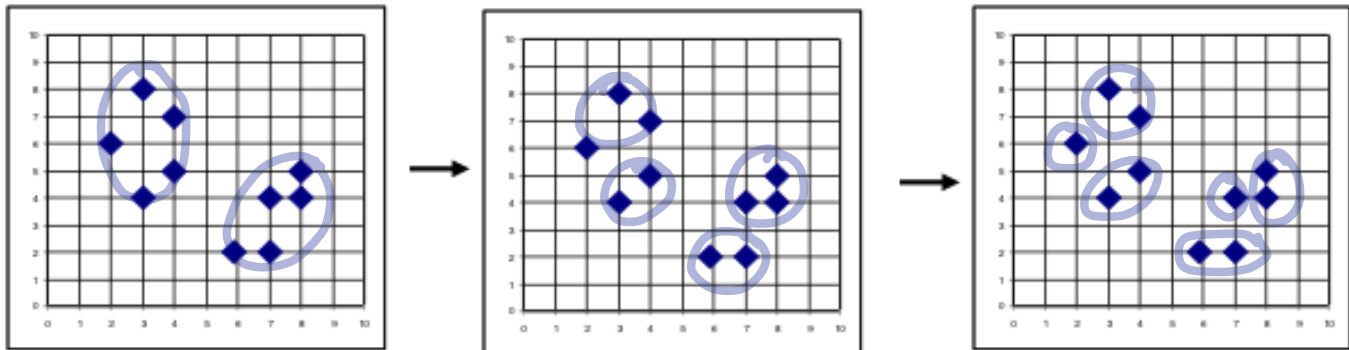
**A clustering of the data objects is obtained by cutting the dendrogram at the desired level, then each connected component forms a cluster.**



# DIANA (Divisive Analysis)

ଭୀମ ଓ ଅକ୍ଷୟ.

- 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





# More on Hierarchical Clustering Methods

- Major weakness of agglomerative clustering methods
  - do not scale well: time complexity of at least  $O(n^2)$ , where  $n$  is the number of total objects
  - can never undo what was done previously 변복 불가능
- Integration of hierarchical with distance-based clustering
  - BIRCH (1996): uses CF-tree and incrementally adjusts the quality of sub-clusters
  - CURE (1998): selects well-scattered points from the cluster and then shrinks them towards the center of the cluster by a specified fraction
  - CHAMELEON (1999): hierarchical clustering using dynamic modeling

이러한 것  
있어  
(개념만)





# Density-Based Clustering Methods

밀도 기반 클러스터링

- Clustering based on density (local cluster criterion), such as density-connected points
- Major features:
  - Discover clusters of arbitrary shape 임의의 모양의 클러스터를 만들 수 있음
  - Handle noise
  - One scan 한번 스캔만으로 가능
  - Need density parameters as termination condition  $\epsilon$  MinPts
- 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)

# Density-Based Clustering: Background

Two parameters:

- **Eps**: Maximum radius of the neighbourhood
- **MinPts**: Minimum number of points in an Eps-neighbourhood of that point

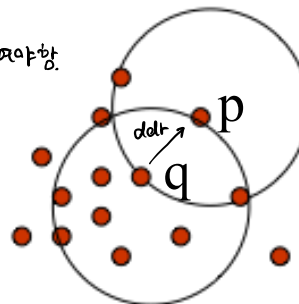
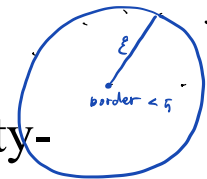
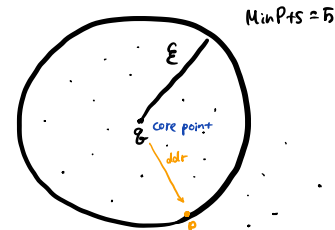
**NEps(p)**:  $\{q \text{ belongs to } D \mid \text{dist}(p, q) \leq \text{Eps}\}$

*p*로부터 일정 반경 안의 이웃들.

**Directly density-reachable**: A point *p* is directly density-reachable from a point *q* wrt. **Eps**, **MinPts** if

- 1) *p* belongs to **NEps(q)** *q*의 이웃이어야 하고
- 2) core point condition: *core point* 여야 함.

$$|NEps(q)| \geq MinPts$$



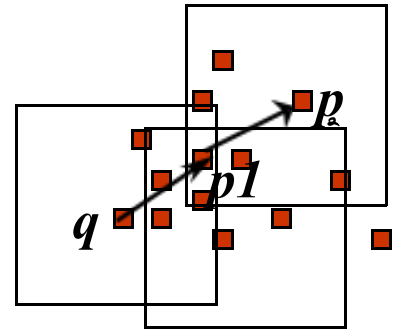
MinPts = 5

Eps = 1 cm

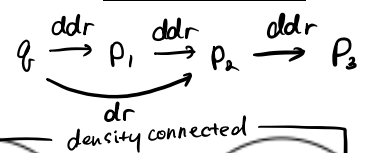
# Density-Based Clustering: Background (II)

## Density-reachable:

A point  $p$  is density-reachable from a point  $q$  wrt.  $Eps$ ,  $MinPts$  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  $p_i$

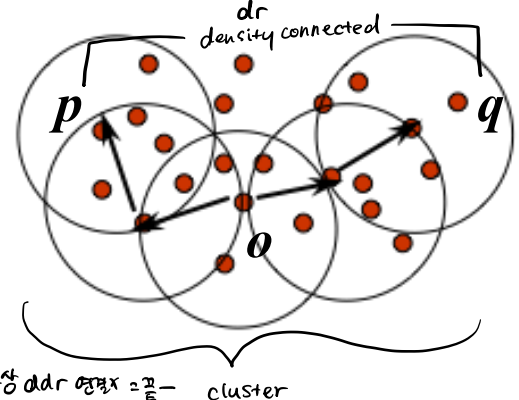


항변 방공식,  
재방문 X



## Density-connected

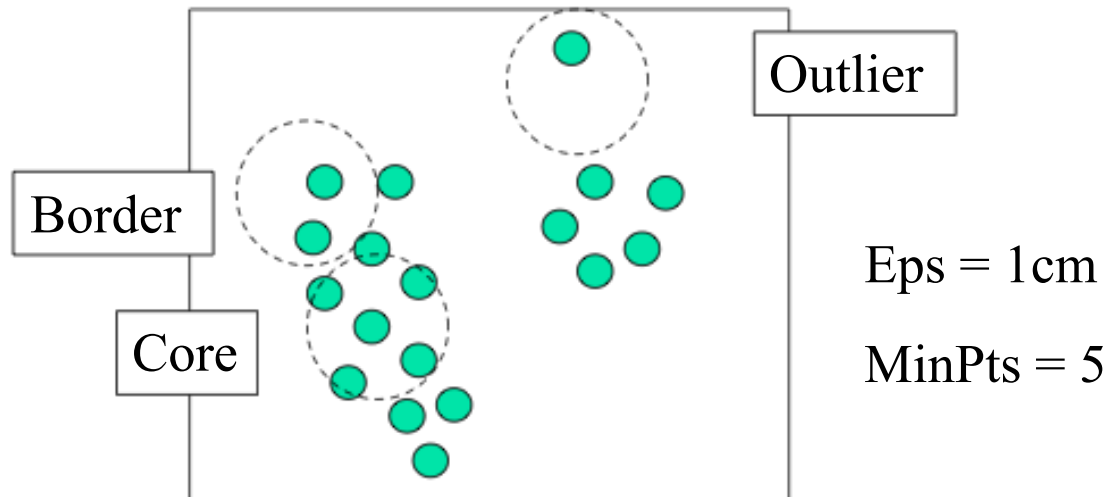
A point  $p$  is density-connected to a point  $q$  wrt.  $Eps$ ,  $MinPts$  if there is a point  $o$  such that both,  $p$  and  $q$  are density-reachable from  $o$  wrt.  $Eps$  and  $MinPts$ .



이러한 ddr 연결은 곧 cluster

# 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





# Grid-Based Clustering Method

---

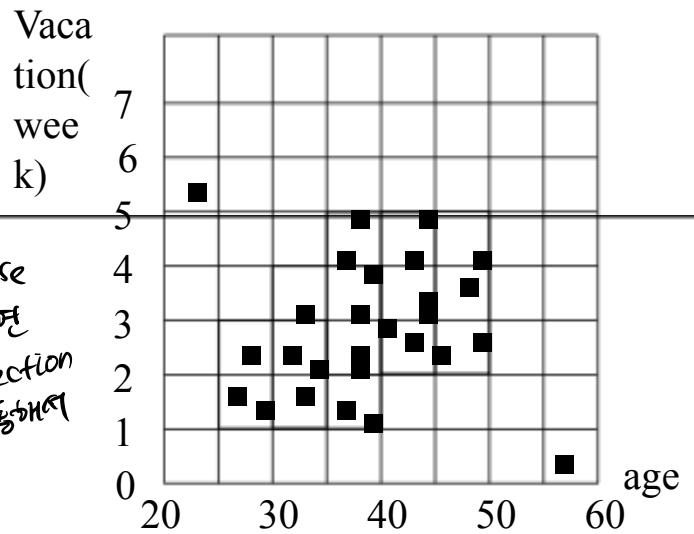
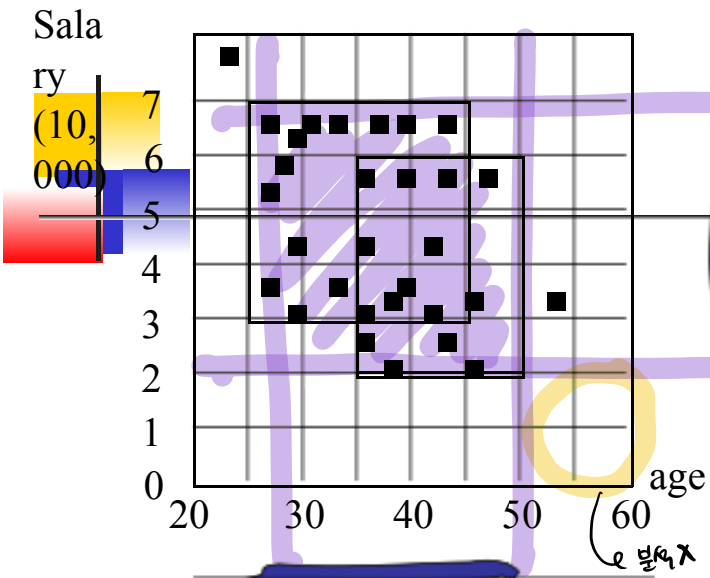
- Using multi-resolution grid data structure
- Several interesting methods
  - **STING** (a S**T**atistical **I**Nformation Grid approach) by Wang, Yang and Muntz (1997)
  - **WaveCluster** by Sheikholeslami, Chatterjee, and Zhang (VLDB'98)
- A multi-resolution clustering approach using wavelet method
- **CLIQUE**: Agrawal, et al. (SIGMOD'98)



# CLIQUE: The Major Steps

---

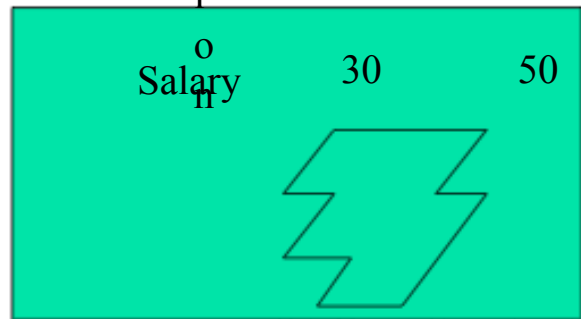
- Partition the data space and find the number of points that lie inside each cell of the partition.
- Identify the subspaces that contain clusters using the Apriori principle
- Identify clusters:
  - Determine dense units in all subspaces of interests
  - Determine connected dense units in all subspaces of interests.
- Generate minimal description for the clusters
  - Determine maximal regions that cover a cluster of connected dense units for each cluster
  - Determination of minimal cover for each cluster



$\square = 3$

탐색공간을 축소.

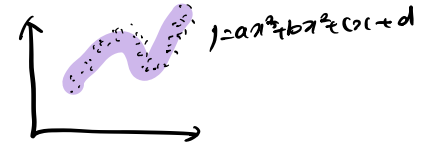
V  
a  
c  
a  
t  
i  
o  
n



age

# Model-Based Clustering Methods

- Attempt to optimize the fit between the data and some mathematical model



- Statistical and AI approach

- Conceptual clustering

- A form of clustering in machine learning

- Produces a classification scheme for a set of unlabeled objects

- Finds characteristic description for each concept (class)

- COBWEB (Fisher'87)

- A popular a simple method of incremental conceptual learning

- Creates a hierarchical clustering in the form of a **classification tree**

- Each node refers to a concept and contains a probabilistic description of that concept



신명양을

PPT 소개하기