



Clustering

강필성

고려대학교 산업경영공학부

Bflysoft & WIGO AI LAB

AGENDA

01 Clustering: Overview

02 K-Means Clustering

03 Hierarchical Clustering

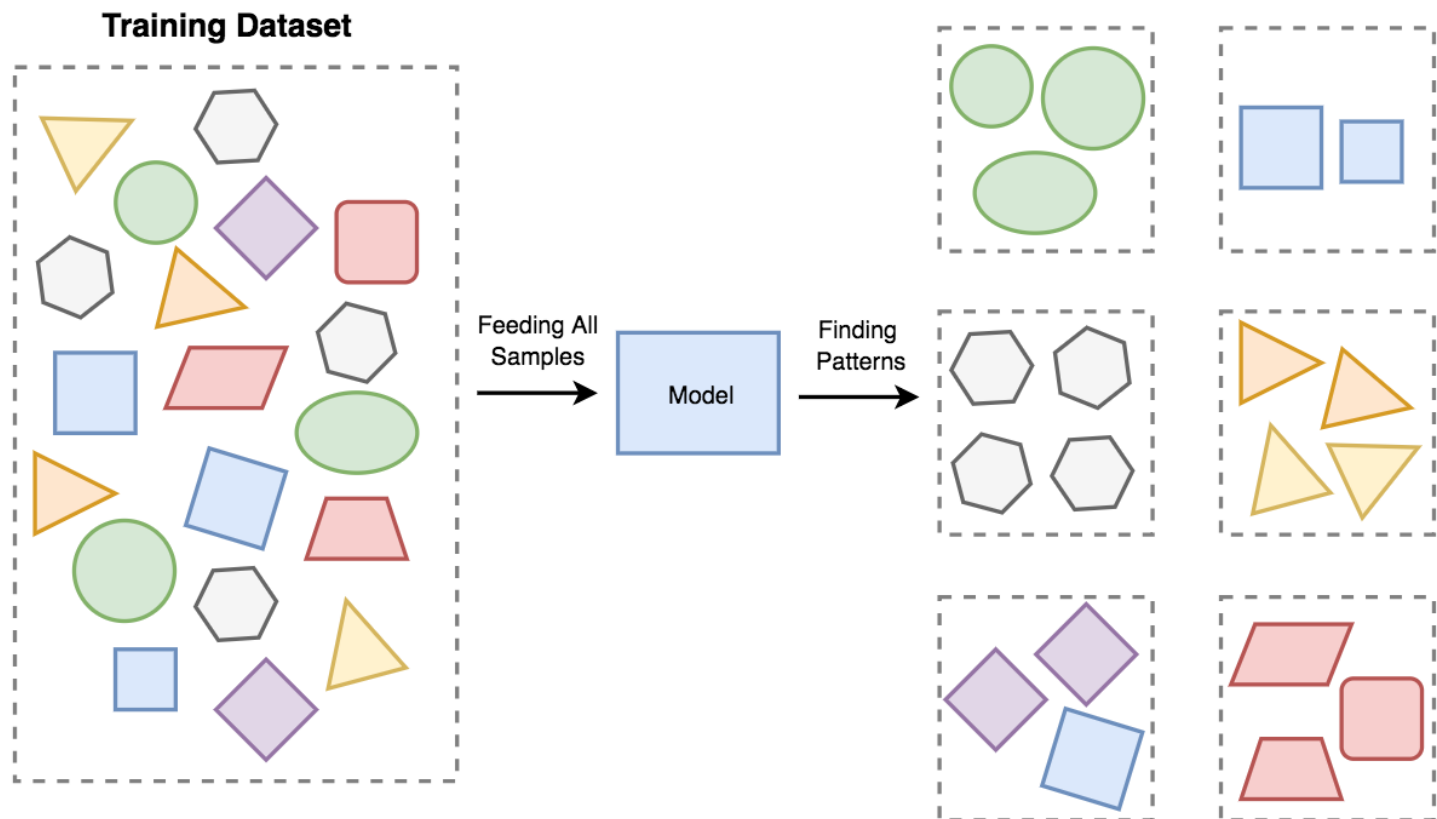
04 Density-based Clustering: DBSCAN

04 R Exercise

Clustering: Overview

- Supervised vs. Unsupervised Learning

- ✓ Supervised: Find a function f that explains the relationship between the input X and the output Y
- ✓ Unsupervised: Explore the features of the input X

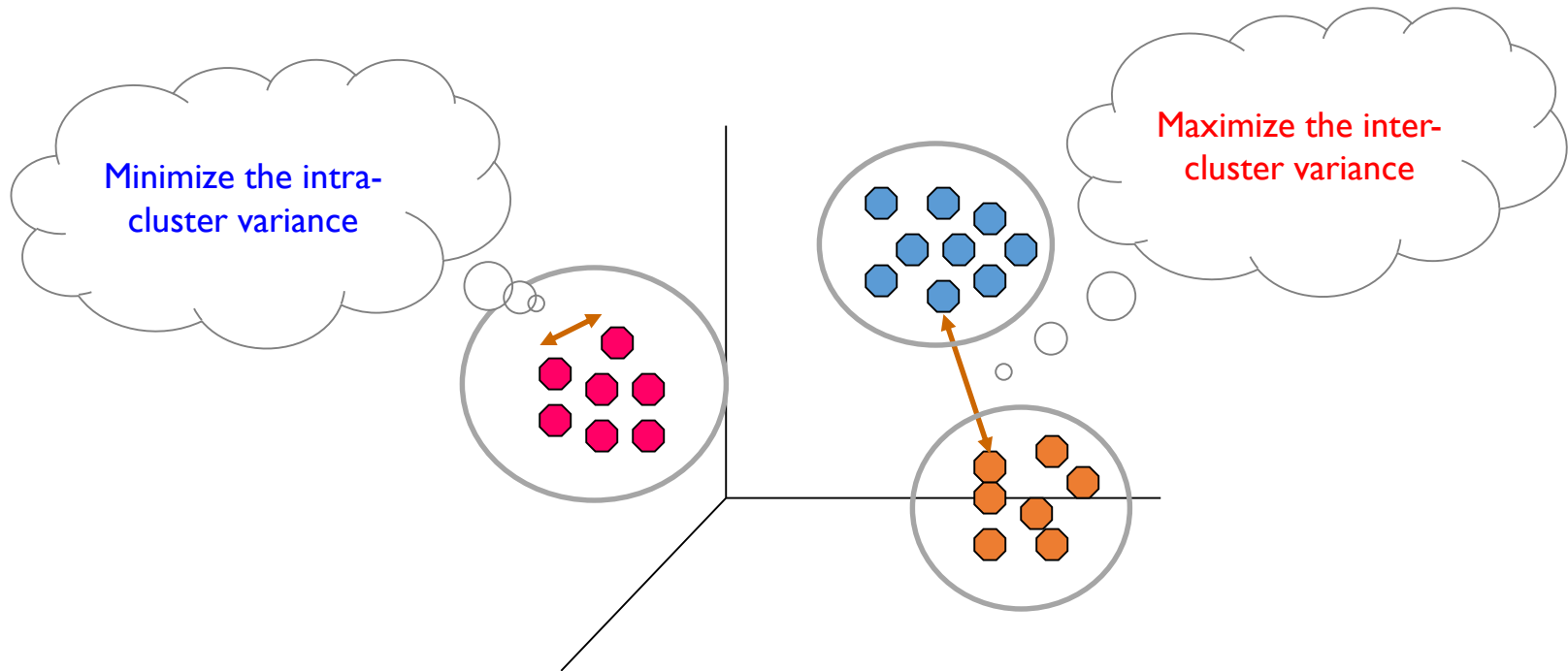


Clustering: Overview

- 군집화(Clustering)

- ✓ 관측치들의 집단을 판별

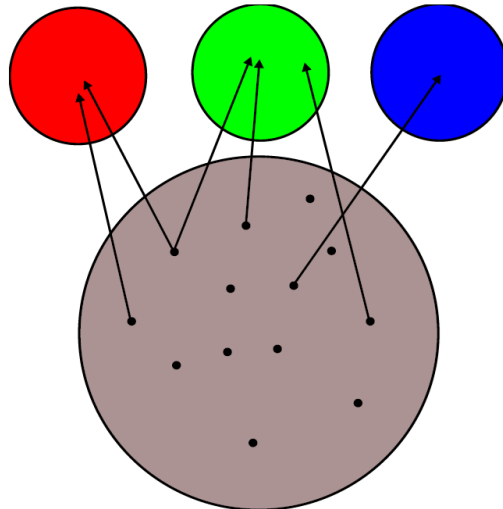
- 동일한 집단에 소속된 관측치들은 서로 유사할수록 좋음
 - 상이한 집단에 소속된 관측치들은 서로 다를수록 좋음



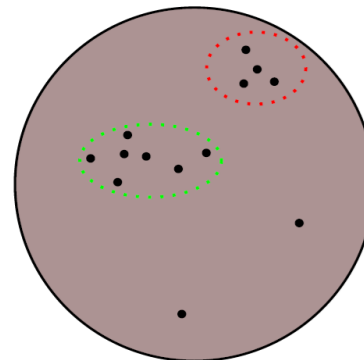
Clustering: Overview

- 분류 (Classification) vs. 군집화(Clustering)

- ✓ **분류(Classification)**: 범주의 수 및 각 개체의 범주 정보를 사전에 알 수 있으며, 개체의 입력 변수 값들로부터 범주 정보를 유추하여 새로운 개체에 대해 가장 적합한 범주로 할당하는 문제 (**supervised learning**)
- ✓ **군집화(Clustering)**: 군집의 수, 속성, 멤버십 등이 사전에 알려져 있지 않으며 최적의 구분을 찾아가는 문제 (**unsupervised learning**)



(a) Classification



(b) Clustering

Clustering: Overview

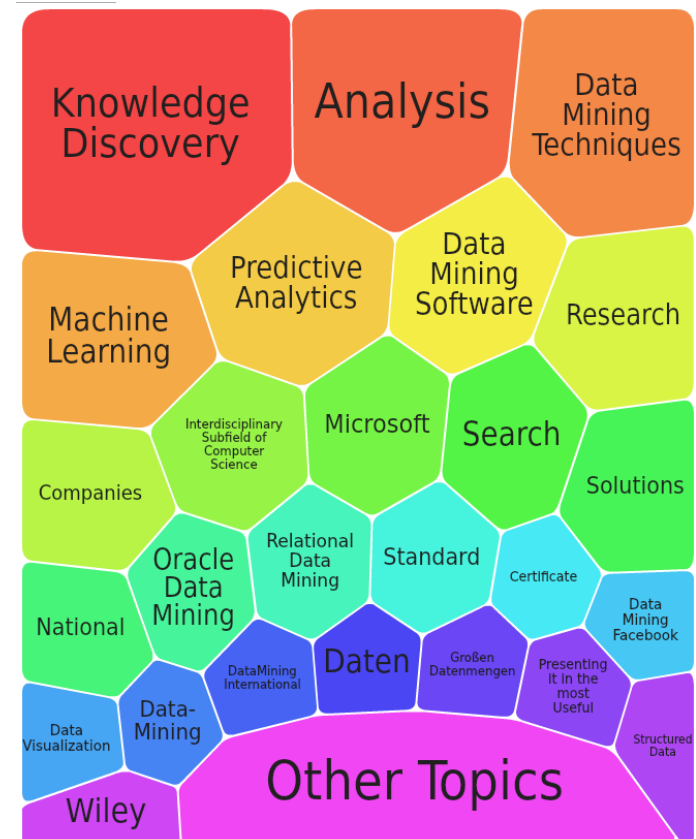
• 군집화 적용 사례

✓ 데이터에 대한 이해

- 웹브라우징 시 유사한 문서들을 표시
- 유사한 기능을 수행하는 유전자/단백질 집합
- 유사한 추세를 나타내는 주식 종목들 등

Query: israel
Documents: 272, Clusters: 15, Average Cluster Size: 15.1 documents

Cluster	Size	Shared Phrases and Sample Document Titles
1 View Results Refine Query Based On This Cluster	16	Society and Culture (56%), Faiths and Practices (56%), Judaism (69%), Spirituality (56%); Religion (56%), organizations (43%) <ul style="list-style-type: none"> ● Ahavat Israel - The Amazing Jewish Website! ● Israel and Judaism ● Judaica Collection
2 View Results Refine Query Based On This Cluster	15	Ministry of Foreign Affairs (33%), Ministry (87%) <ul style="list-style-type: none"> ● Publications and Data of the BANK OF ISRAEL ● Consulate General of Israel to the Mid-Atlantic Region ● The Friends of Israel Gospel Ministry
3 View Results Refine Query Based On This Cluster	11	Israel Tourism (36%), Comprehensive Israel (36%), Tourism (64%) <ul style="list-style-type: none"> ● Interactive Israel tourism guide - Jerusalem ● Ambassade d'Israel ● Travel to Israel Opportunities
4 View Results Refine Query Based On This Cluster	7	Middle East (57%), History (57%); WAR (42%), Region (42%), Complete (42%), Listing (42%), country (42%) <ul style="list-style-type: none"> ● Israel at Fifty: Our Introduction to The Six Day War ● Machal - Volunteers in the Israel's War of Independence ● HISTORY: The State of Israel
5 View Results Refine Query Based On This Cluster	22	Economy (68%), Companies (55%), Travel (55%) <ul style="list-style-type: none"> ● Israel Hotel Association ● Israel Association of Electronics Industries ● Focus Capital Group - Israel



Clustering: Overview

- 군집화 적용 사례

- ✓ 전략 수립

- 군집화를 이용한 자산 관리(asset management)
 - 개별 종목의 과거 6개월 수익률 및 변동성을 변수로 사용하여 계층적 군집화 수행
 - 네 개의 군집 중에서 maximum performance와 minimum volatility를 나타내는 군집에 속한 종목들로 포트폴리오 구성



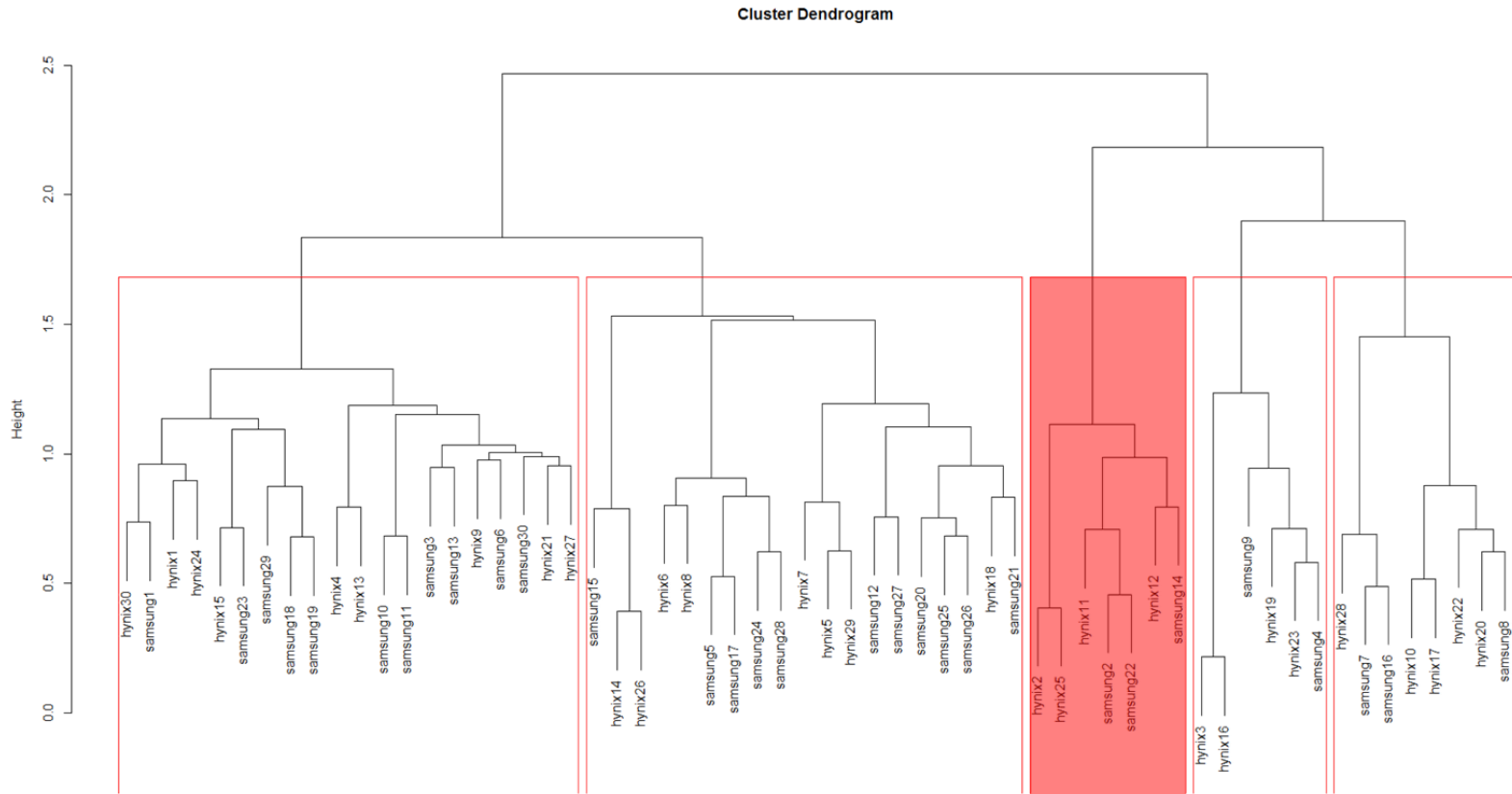
<https://quantdare.com/hierarchical-clustering/>

Clustering: Overview

- 군집화 적용 사례

- ✓ 전략 수립

- 경쟁사와의 특허 문서 분석을 통한 장단점 파악



as.dist(1 - cosine(normMat))
hclust("ward.D")

Clustering: Overview

• 군집화 적용 사례

✓ 전략 수립

▪ 경쟁사와의 특허 문서 분석을 통한 장단점 파악

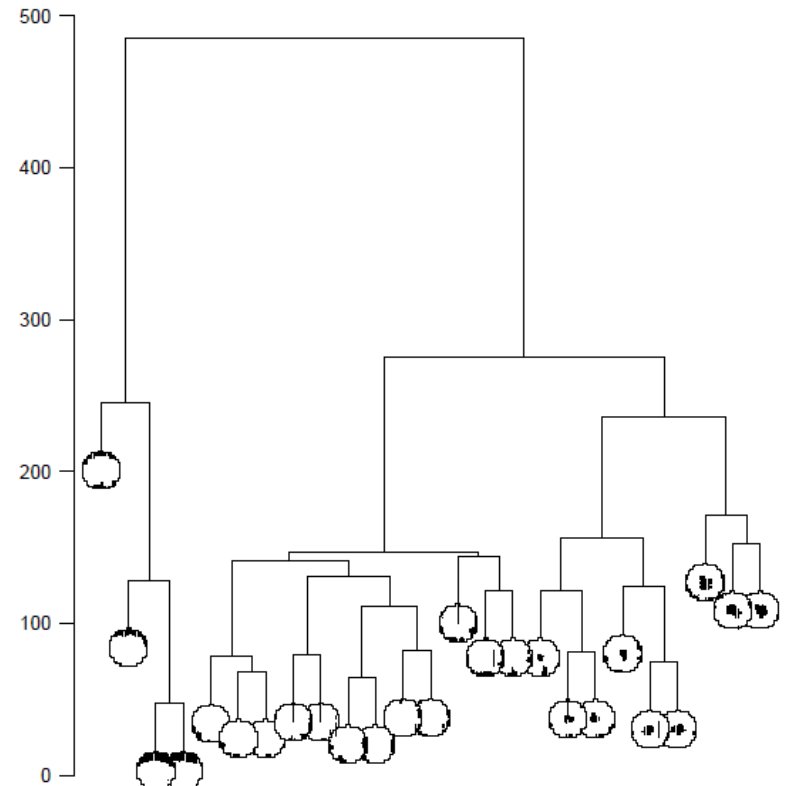
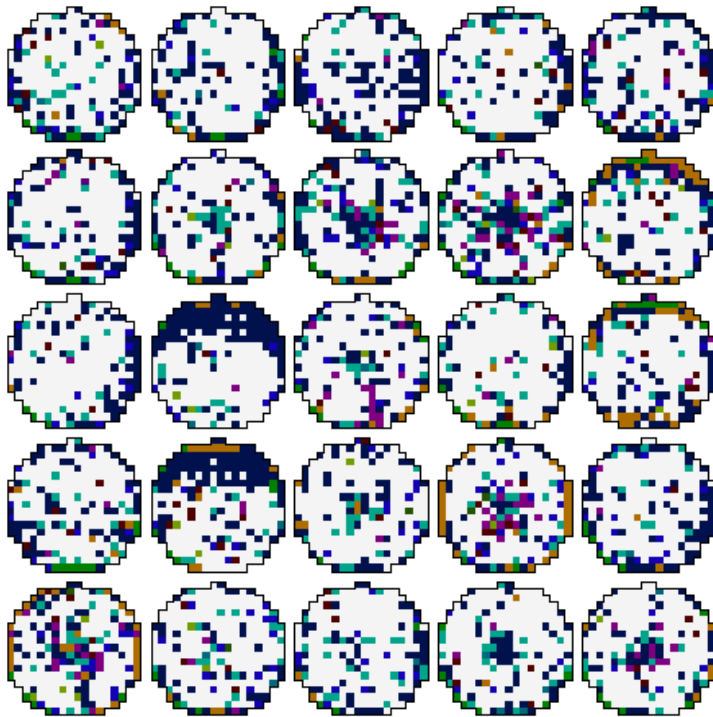
1	회사	일련번호	특허명	조록
2	SK하이닉스	2	멀티 레글레이터 회로 및 이를 구비한 집적회로	본 기술에 따른 레글레이터 회로는, 입력전압을 일정한 전압 레벨로 레글레이팅하여 출력하도록 구성된 레글레이터 및 복수개의 전압 생성 코드 들에 의해 결정되는 내부 저항값들에 따라 상기 레글레이터의 출력 전압을 분배한 분배전압들을 각각 출력하도록 구성된 복수개의 전압 분배회로를 포함한다.
3	SK하이닉스	11	내부 전압 생성 회로 및 그의 동작 방법	펄핑 동작을 통해 내부 전압을 생성하는 내부 전압 생성 회로에 관한 것으로, 다수의 펄핑부를 포함하며, 목표 전압 레벨에 대응하는 최종 펄핑 전압을 생성하기 위한 펄핑 전압 생성부, 및 상기 목표 전압 레벨에 대응하여 상기 다수의 펄핑부의 활성화 개수를 제어하기 위한 활성화 제어부를 구비하는 내부 전압 생성 회로가 제공된다.
4	SK하이닉스	12	자기 메모리 장치를 위한 라이트 드라이버 회로 및 자기 메모리 장치	비트라인과 소스라인 간에 접속되며, 비트라인 방향으로 인접하는 한 쌍의 자기 메모리 셀이 소스라인을 공유하는 복수의 자기 메모리 셀로 이루어진 메모리 셀 어레이를 포함하는 자기 메모리 장치를 위한 라이트 드라이버 회로로서, 정의 기록전압 공급단자와 부의 기록전압 공급단자 간에 접속되어, 라이트 인에이블 신호 및 데이터 신호에 따라 정의 기록전압 또는 부의 기록전압에 의한 전류를 비트라인에 선택적으로 공급하는 스위칭부를 포함하는 자기 메모리 장치를 제공한다.
5	SK하이닉스	25	전압 레글레이터 및 전압 레글레이팅 방법	전압 레글레이터는 출력전압을 전압 출력단으로 출력하는 전압 출력부와, 제1 제어코드의 제어에 따라 분배 저항값을 조절하는 제1 저항분배 스테이지와, 제1 저항분배 스테이지에서 결정된 분배 저항값을 제2 제어코드의 제어에 따라 조절하는 제2 저항분배 스테이지를 포함하며, 전압 출력단을 통해서 출력되는 출력전압의 전압레벨은 제1 및 제2 저항분배 스테이지를 통해서 결정된 상기 분배 저항값과, 기준저항의 저항값 비율에 따라 조절되는 것을 특징으로 한다.
6	삼성전자	2	전압 공급 장치 및 그것을 포함한 불휘발성 메모리 장치	본 발명에 따른 전압 공급 장치는 전원 전압을 승압하고, 상기 승압된 전압을 출력 라인으로 제공하기 위한 전하 펌프 및 상기 출력 라인의 전압 레벨을 목표 전압 레벨로 유지하기 위한 전압 제어 회로를 포함한다. 본 발명에 따른 상기 전압 제어 회로는 웰 상에 형성된 제 1 영역 및 제 2 영역을 포함하고, 상기 제 1 영역 및 제 2 영역 사이의 리치 스루(reach through)를 이용하여 상기 출력 라인의 전압 레벨을 제어하기 위한 리치 스루 소자를 포함한다.
7	삼성전자	14	파워 공급 회로 및 이를 구비하는 상 변화 메모리 장치	파워 공급 회로 및 이를 구비하는 상 변화 메모리 장치가 개시된다. 본 발명의 제 1 실시예에 따른 반도체 메모리 장치는 파워 공급 회로, 스위치들 및 선택기들을 구비한다. 파워 공급 회로는 상기 블록들의 메모리 셀들에 사용되는 제 1 전압 및 제 2 전압을 생성한다. 스위치들은 상기 파워 공급 회로와 상기 제 1 전압이 전달되는 제 1 라인 및 상기 제 2 전압이 전달되는 제 2 라인으로 연결되고, 제어 신호에 응답하여 상기 제 1 전압 및 제 2 전압 중 하나를 대응되는 블록으로 인가한다. 선택기들은 블록 선택 신호 및 디스차이지 성공 신호에 응답하여, 상기 제어 신호를 생성한다. 본 발명에 따른 파워 공급 회로 및 이를 구비하는 상 변화 메모리 장치는 셀 블록마다 별도의 파워 스위치를 구비함으로써 파워 공급 회로의 동작 시간 및 동작 전류를 감소시킬 수 있다. 또한, 기입 전압을 디스차이지한 후 다른 레벨의 전압을 공급함으로써, 상 변화 메모리 장치의 오작동이 방지될 수 있다.
8	삼성전자	22	전압 안정화 장치 및 그것을 포함하는 반도체 장치 및 전압 생성 방법	본 발명은 전압 안정화 장치 및 그것을 이용하는 반도체 장치에 관한 것이다. 본 발명의 기술적 사상의 실시예에 따른 전압 안정화 장치는 제 1 전압을 생성하는 제 1 레글레이터 및 상기 제 1 전압보다 낮은 제 2 전압을 생성하는 제 2 레글레이터를 포함하고, 상기 제 2 레글레이터는 상기 제 1 전압의 레벨과 미리 정해진 기준 전압의 레벨의 비교 결과에 기초하여 상기 제 1 전압 또는 상기 제 1 전압보다 높은 제 2 전압을 선택적으로 이용하여 상기 제 2 전압을 생성한다. 본 발명의 기술적 사상의 실시예에 따르면 제 1의 전압 > 제 2의 전압의 관계를 유지하면서, 동시에 제 2의 전압을 고속으로 전위 변환시킬 수 있다.

Clustering: Overview

- 군집화 적용 사례

- ✓ 웨이퍼 Fail bit map 군집화: 불량 패턴의 군집화

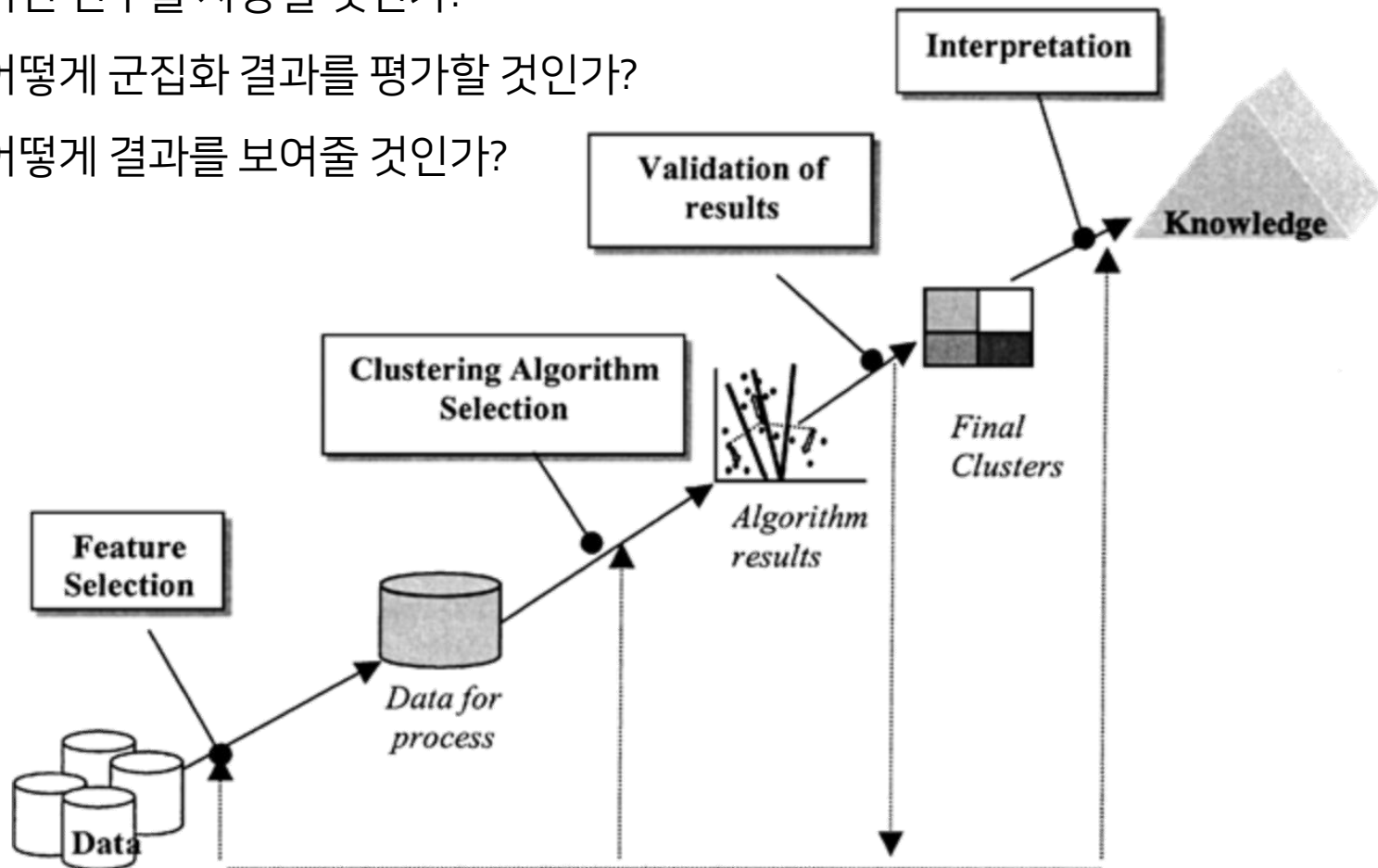
Sample lot exhibiting spatial patterning



Clustering: Overview

- 일반적인 군집화 수행 절차

- ✓ 어떤 변수를 사용할 것인가?
- ✓ 어떻게 군집화 결과를 평가할 것인가?
- ✓ 어떻게 결과를 보여줄 것인가?



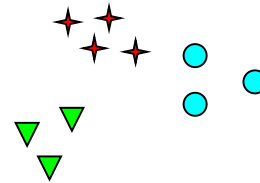
Clustering: Issues

- 군집화 수행 시 고려 사항

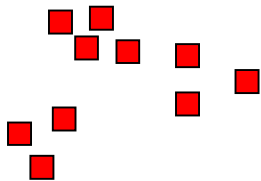
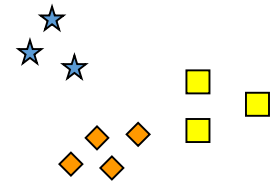
✓ 최적의 군집 수 결정



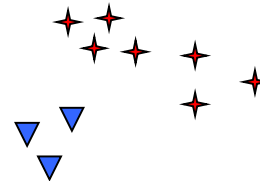
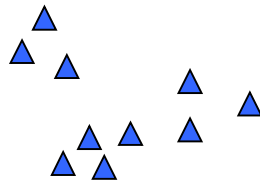
How many clusters?



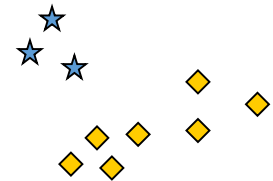
Six Clusters



Two Clusters



Four Clusters

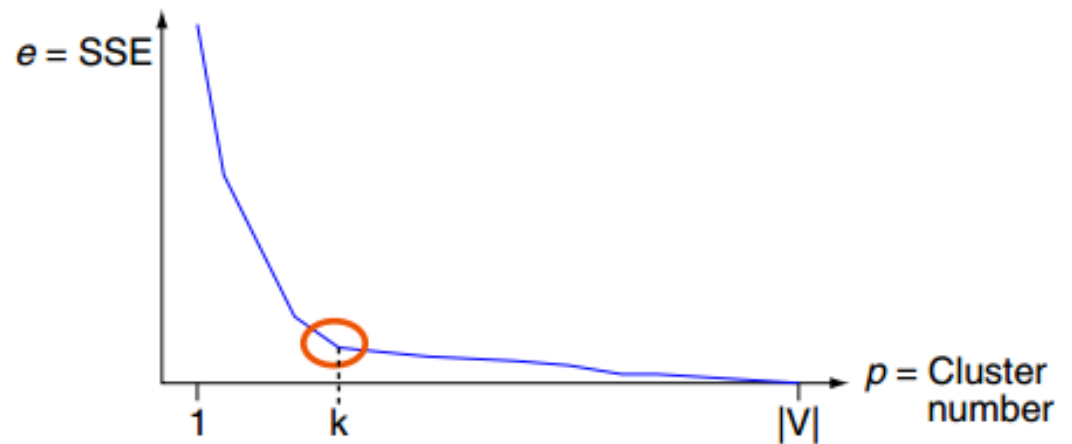
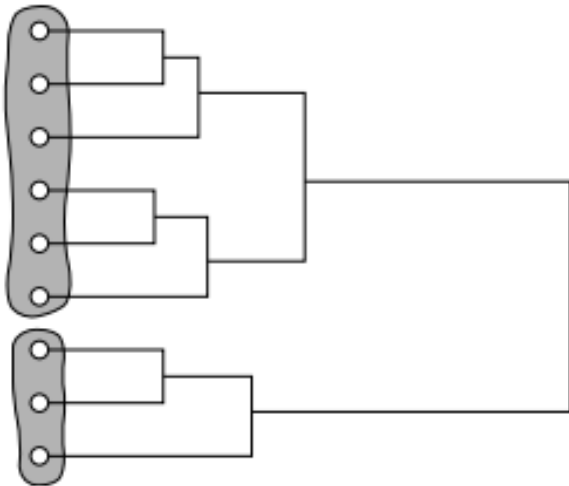


Clustering: Issues

- 군집화 수행 시 고려 사항

- ✓ 최적의 군집 수 결정

- 다양한 군집 수에 대해 성능 평가 지표를 도출하여 최적의 군집 수 선택
 - Elbow point에서 최적 군집 수가 결정되는 경우가 일반적임



Clustering: Issues

- 군집화 수행 시 고려 사항
 - ✓ 군집화 결과를 어떻게 평가할 것인가?
 - ✓ 분류/회귀 알고리즘처럼 모든 상황에서 적용가능한 Global Performance Measure 부재
- 군집화 평가 지표는 다음과 같이 세 가지 카테고리로 구분할 수 있음
 - ✓ External: 정답 레이블과의 비교를 통해 성능 평가 (현실적으로 불가능)
 - ✓ Internal: “군집이 얼마나 컴팩트한가”에 보다 초점을 둠
 - ✓ Relative: “군집이 얼마나 컴팩트한가”와 “군집끼리 얼마나 다른가”를 동시에 고려하고자 함

Clustering: Issues

- 군집화 수행 시 고려 사항

- ✓ 군집화 결과를 어떻게 평가할 것인가?

- ✓ 분류/회귀 알고리즘처럼 모든 상황에서 적용가능한 Global Performance Measure 부재

External



- ☐ Rand Statistic
- ☐ Jaccard Coefficient
- ☐ Folks and Mallows index
- ☐ (Normalized) Hurbert Γ statistic

Internal



- ☐ Cophenetic Correlation Coefficient
- ☐ Sum of Squared error (SSE)
- ☐ Cohesion and separation

Relative



- ☐ Dunn family of indices
- ☐ Davies-Bouldin (DB) index
- ☐ Semi-partial R-squared
- ☐ SD validity index
- ☐ Silhouette

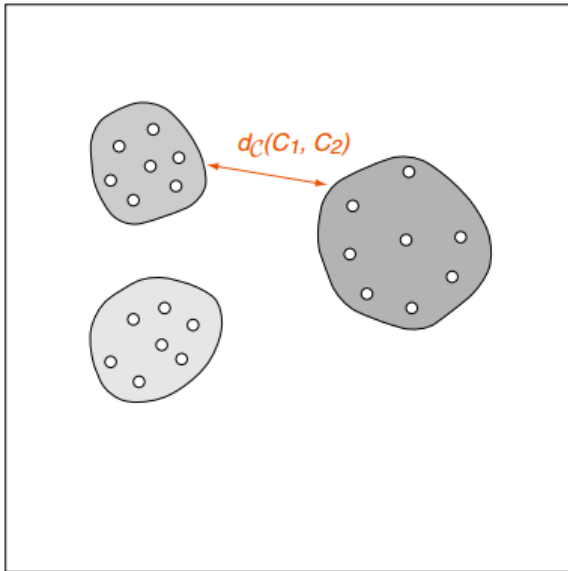
Clustering: Issues

- 군집화 평가를 위해 필요한 세 가지 지표 정의

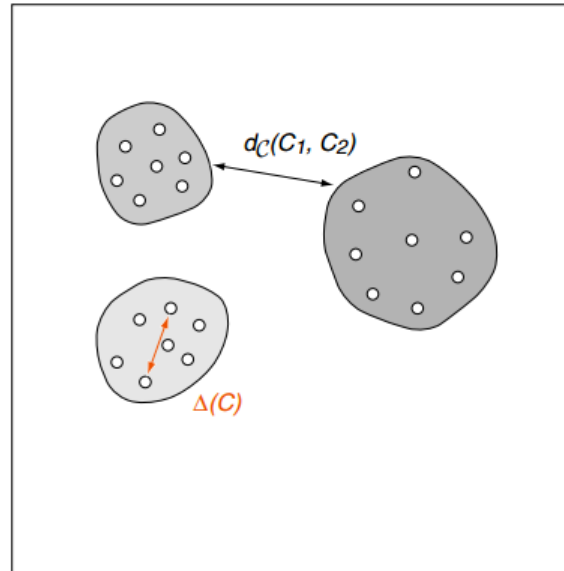
- ✓ 군집화가 잘 되어 있다면

- (1)번 지표의 값은 크고 (2)번과 (3)번 지표의 값은 상대적으로 작게 나타날 것임

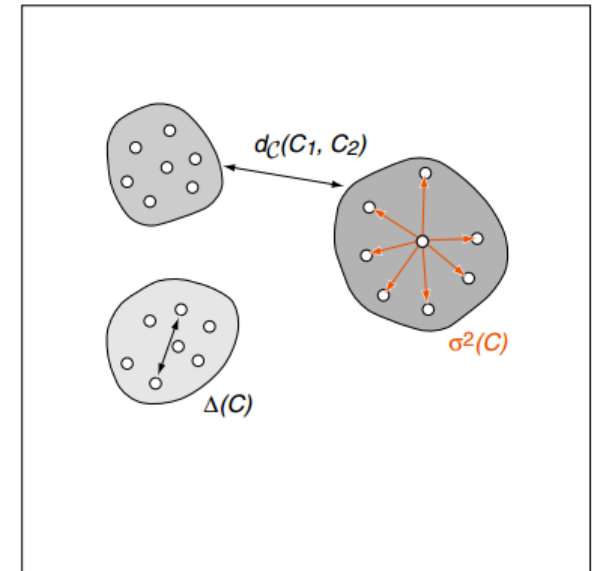
(1) Distance between two clusters



(2) Diameter of a cluster



(3) Scatter within a cluster (SSE)



Clustering: Issues

- 군집화 평가 지표 1: Dunn Index

- ✓ Dunn Index는 군집 내 거리((1)번 지표)중 가장 작은 값을 분자로, 군집의 지름((2)번 지표) 중 가장 큰 값을 분모로 정의함
- ✓ Dunn Index는 클수록 우수한 군집화 결과라고 할 수 있음

$$I(C) = \frac{\min_{i \neq j} \{d_c(C_i, C_j)\}}{\max_{1 \leq l \leq k} \{\Delta(C_l)\}},$$

$I(C) \rightarrow \max$

$$I(C) = \frac{\min_{i \neq j} \{d_c(C_i, C_j)\}}{\max_{1 \leq l \leq k} \{\Delta(C_l)\}},$$

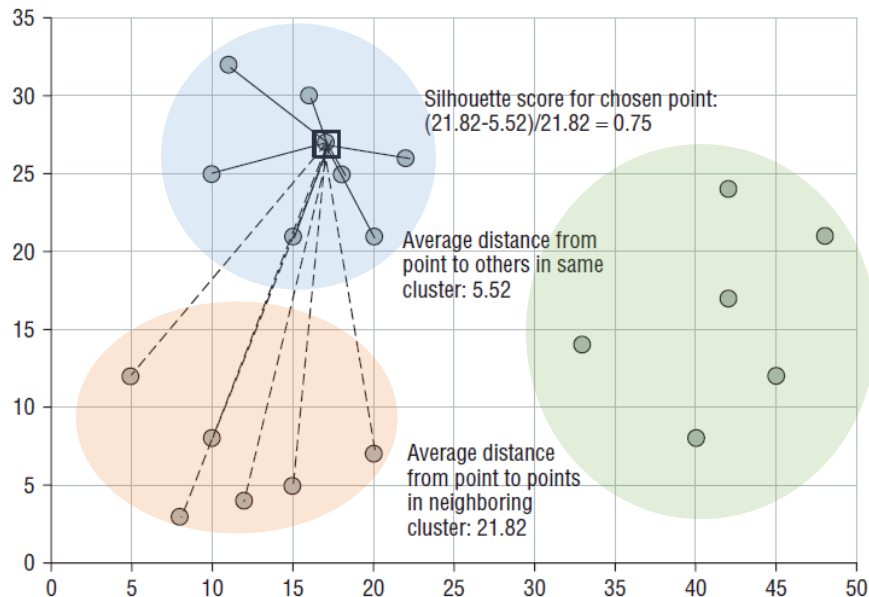
$I(C) \rightarrow \max$

Clustering: Issues

- 군집화 평가 지표 2: Silhouette

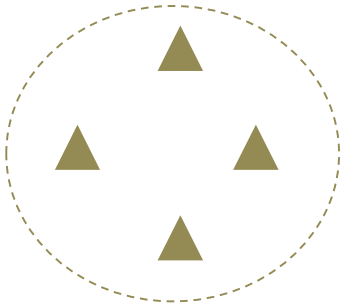
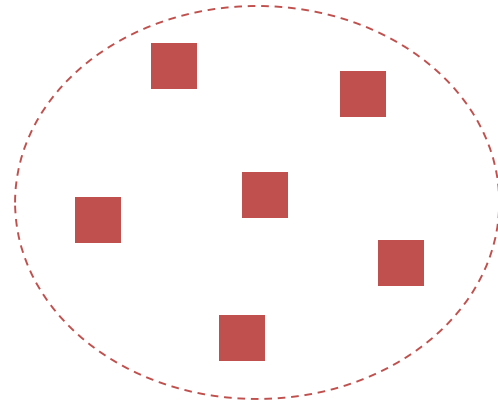
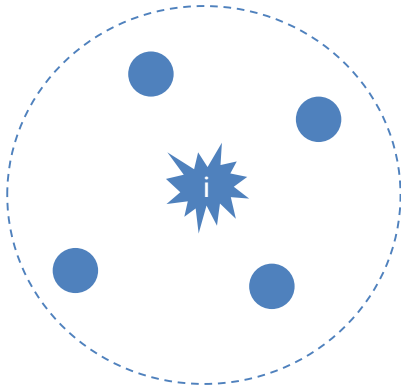
- ✓ $a(i)$: 개체 i 로부터 같은 군집 내에 있는 모든 다른 개체들 사이의 평균 거리
- ✓ $b(i)$: 개체 i 로부터 다른 군집 내에 있는 개체들 사이의 평균 거리 중 가장 작은 값

$$s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}} \quad s(i) = \begin{cases} 1 - a(i)/b(i), & \text{if } a(i) < b(i) \\ 0, & \text{if } a(i) = b(i) \\ b(i)/a(i) - 1, & \text{if } a(i) > b(i) \end{cases}$$



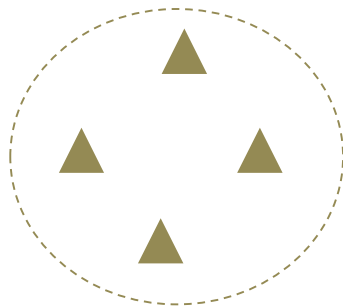
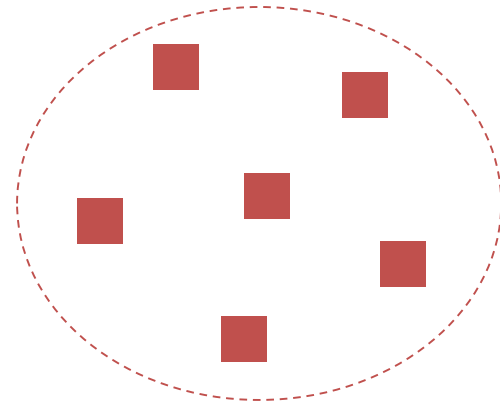
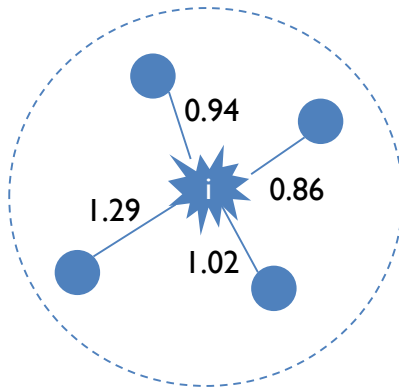
Clustering: Issues

- 군집화 평가 지표 2: Silhouette 계산 예시



Clustering: Issues

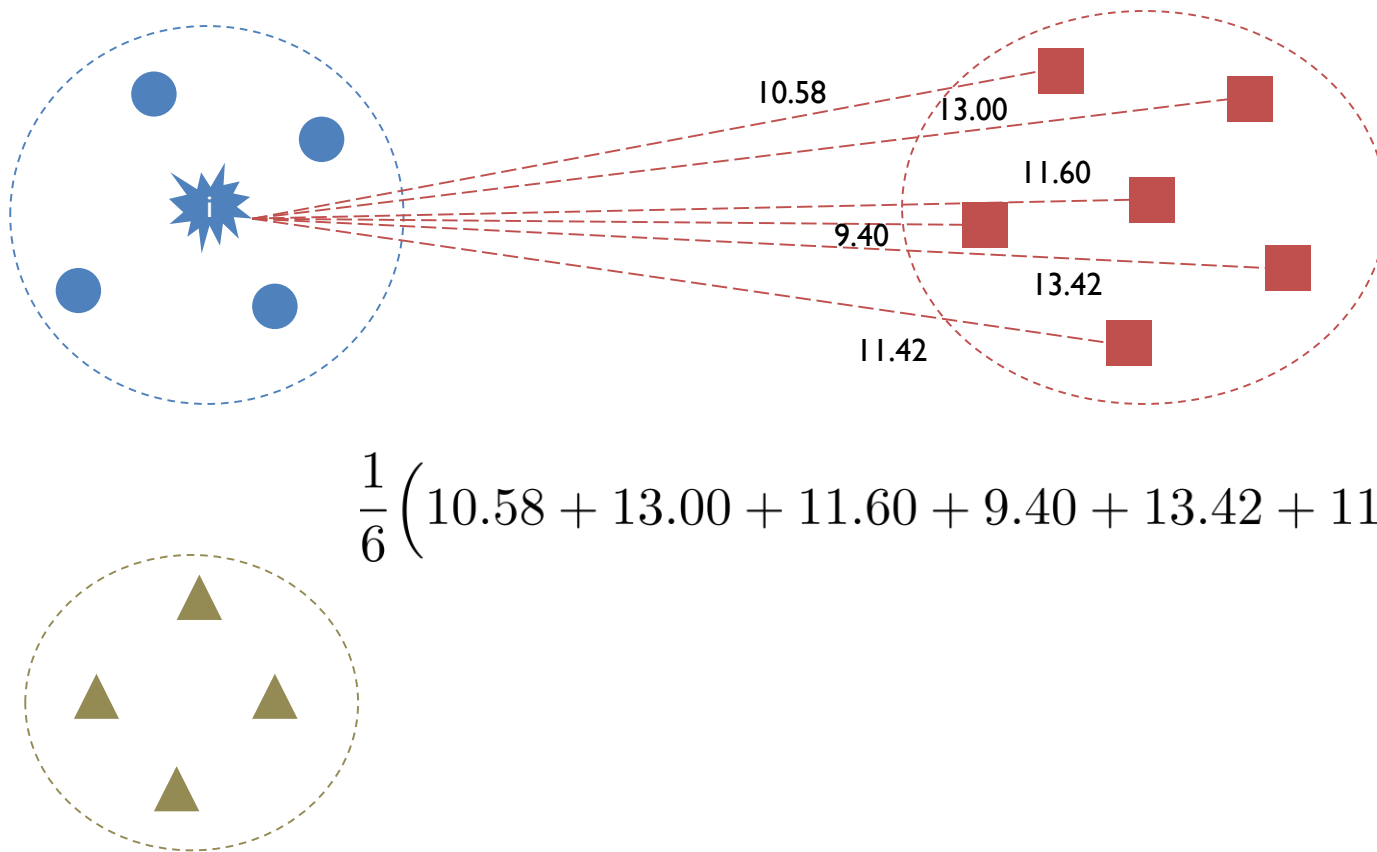
- 군집화 평가 지표 2: Silhouette 계산 예시



$$a(i) = \frac{1}{4} (0.94 + 0.86 + 1.02 + 1.29) = 1.03$$

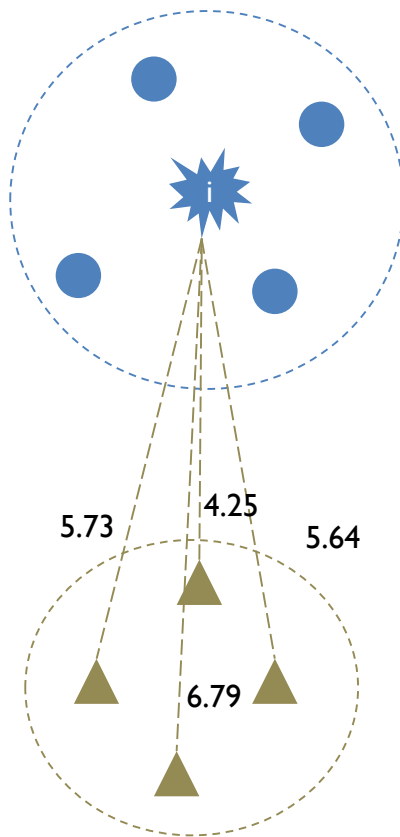
Clustering: Issues

- 군집화 평가 지표 2: Silhouette 계산 예시

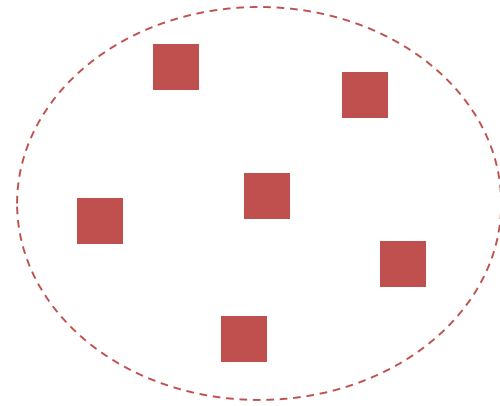


Clustering: Issues

- 군집화 평가 지표 2: Silhouette 계산 예시



$$\frac{1}{4} (5.73 + 6.79 + 4.25 + 5.64) = 5.60$$



$$b(i) = \min(11.57, 5.60) = 5.60$$

$$s(i) = \frac{5.60 - 1.03}{\max(1.03, 5.60)}$$

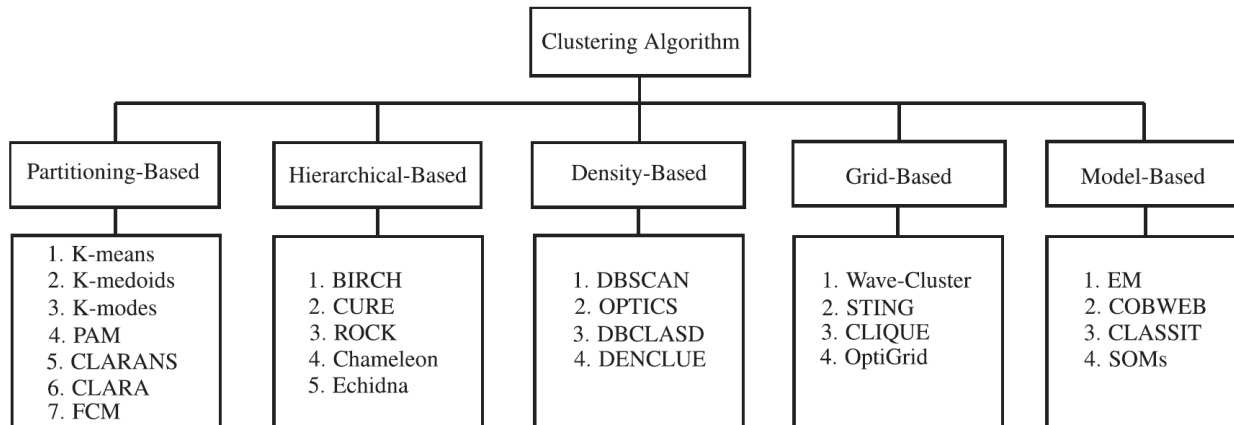
$$= \frac{4.57}{5.60} = 0.82$$

Clustering: Types

- Hard Clustering vs. Soft Clustering

- ✓ Hard Clustering (Crisp Clustering)

- 서로 겹치지 않는(non-overlapping) 군집 생성
 - 각 개체는 오직 하나의 군집으로만 할당됨



- ✓ Soft Clustering (Fuzzy Clustering)

- 겹치는 군집을 생성하는 것도 가능함
 - 한 개체는 여러 개의 군집에 확률적인 할당이 될 수 있음

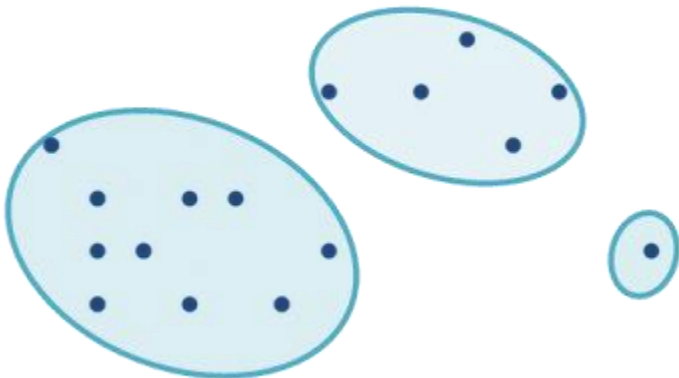
Clustering: Algorithms

• 군집화 알고리즘의 종류

✓ 분리형 군집화

- 전체 데이터의 영역을 특정 기준에 의해 동시에 구분
- 각 개체들은 사전에 정의된 군집 수 중 하나에 속하는 결과를 도출함

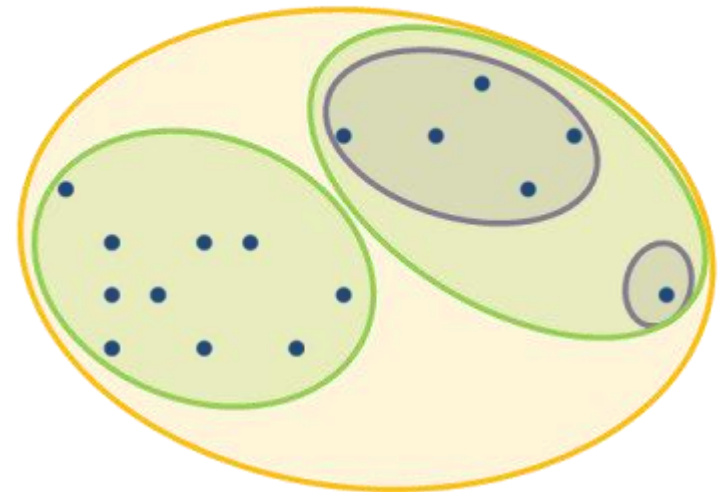
Partitional Clustering



✓ 계층적 군집화

- 개체들을 가까운 집단부터 차근차근 묶어가는 방식
- 군집화 결과 뿐만 아니라 유사한 개체들이 결합되는 절차(dendrogram)도 생성

Hierarchical Clustering



AGENDA

01 Clustering: Overview

02 **K-Means Clustering**

03 Hierarchical Clustering

04 Density-based Clustering: DBSCAN

K-Means Clustering

- K-평균 군집화

- ✓ 대표적인 분리형 군집화 알고리즘

- 각 군집은 하나의 중심(centroid)을 가짐
- 각 개체는 가장 가까운 중심에 할당되며, 같은 중심에 할당된 개체들이 모여 하나의 군집을 생성
- 사전에 군집의 수 K가 정해져야 알고리즘을 실행할 수 있음

$$\mathbf{X} = C_1 \cup C_2 \dots \cup C_K, \quad C_i \cap C_j = \phi, \quad i \neq j$$

$$\arg \min_{\mathbf{C}} \sum_{i=1}^K \sum_{\mathbf{x}_j \in C_i} \|\mathbf{x}_j - \mathbf{c}_i\|^2$$

K-Means Clustering

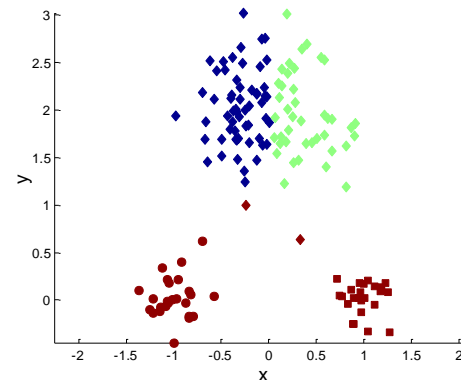
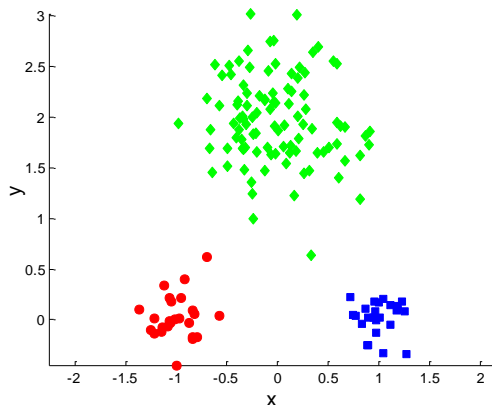
- K-평균 군집화 (K-Means Clustering) 수행 절차

- ✓ 1단계: K개의 초기 군집 중심(initial centroid) 설정

- ✓ 2단계: 다음 절차를 반복

- 모든 개체를 가장 가까운 군집 중심에 할당하여 군집 구성
- 할당된 개체들을 이용하여 군집 중심을 재설정
- 종료 조건: 모든 군집 중심의 위치가 변하지 않고, 모든 개체의 군집 할당 결과에 변화가 없을 때 알고리즘 종료

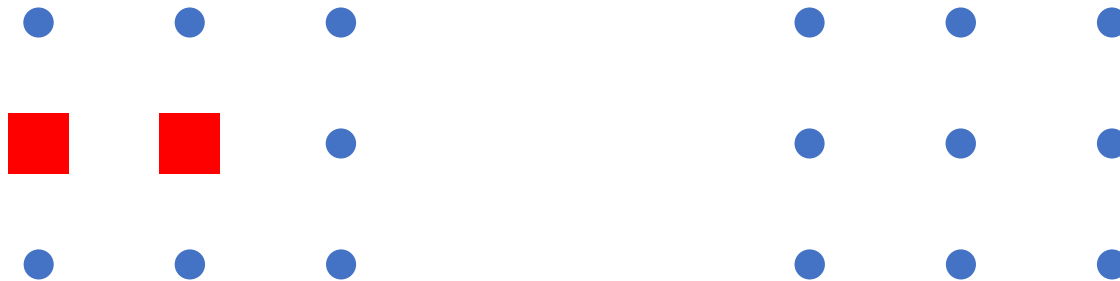
- ✓ Note: 초기 중심은 종종 무작위로 선택되며 따라서 군집화의 결과가 초기 중심 설정에 따라 다르게 나타나는 경우가 발생할 수도 있음



K-Means Clustering

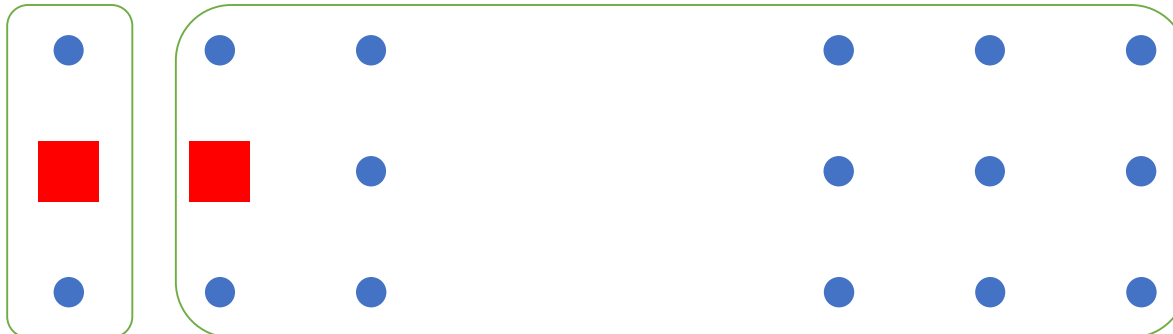
- K-평균 군집화 (K-Means Clustering) 수행 절차 예시

- ✓ 1단계: K개의 초기 군집 중심(initial centroid) 설정



- ✓ 2-1단계(1회차): 모든 개체를 가장 가까운 중심에 할당

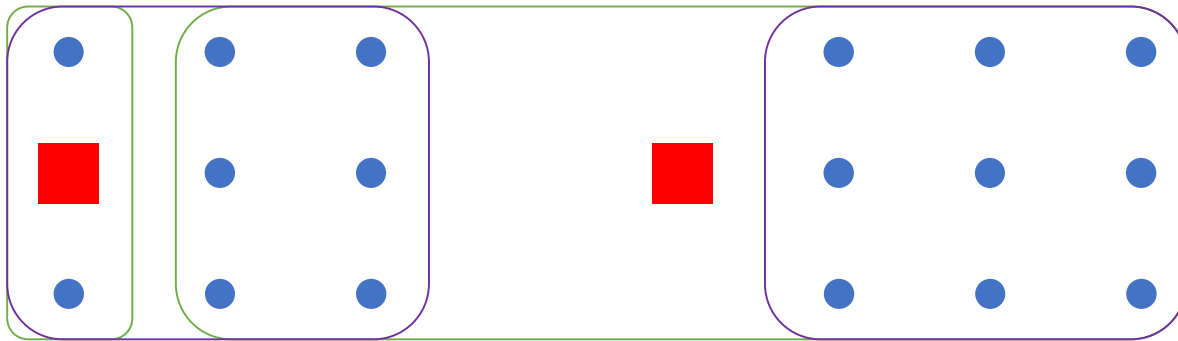
- ✓ 2-2단계(2회차): 할당된 개체들을 이용하여 군집 중심 재설정



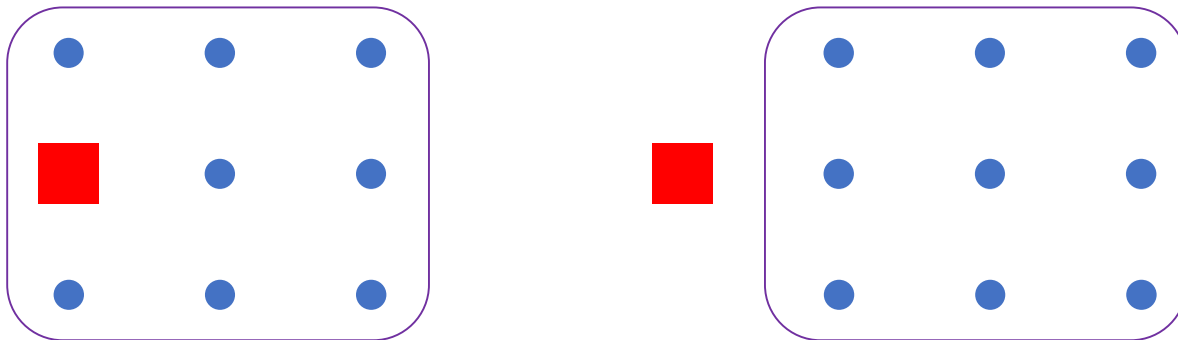
K-Means Clustering

- K-평균 군집화 (K-Means Clustering) 수행 절차 예시

- ✓ 2-1단계(2회차): 모든 개체를 가장 가까운 중심에 할당



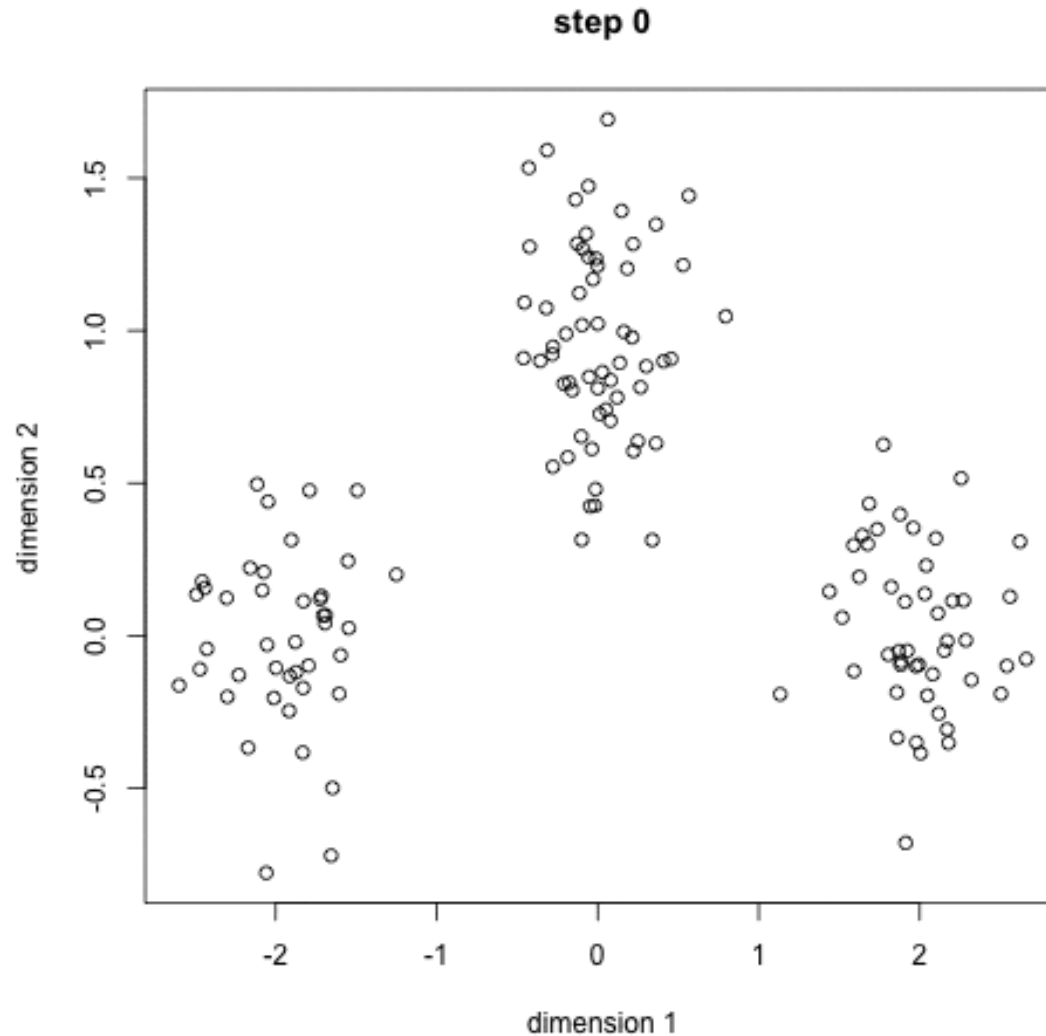
- ✓ 2-2단계(2회차): 할당된 개체들을 이용하여 군집 중심 재설정



- ✓ 군집 중심과 개체 할당에 변화가 없으므로 알고리즘 종료

K-Means Clustering

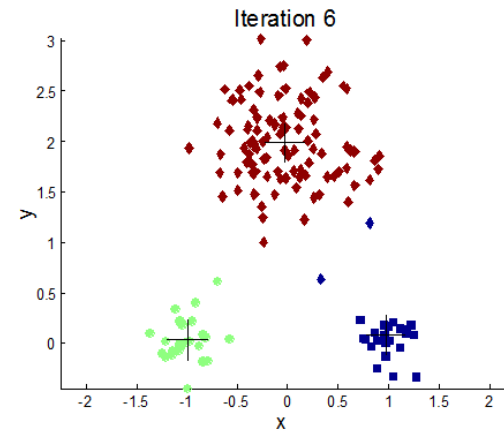
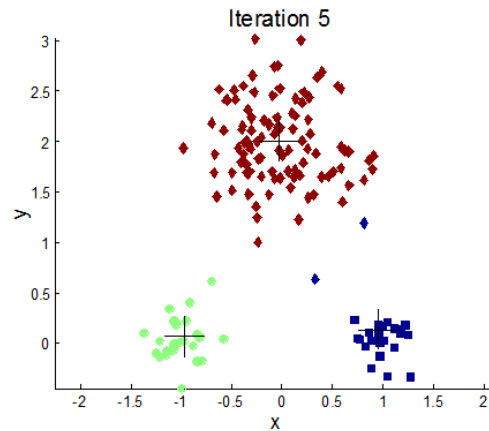
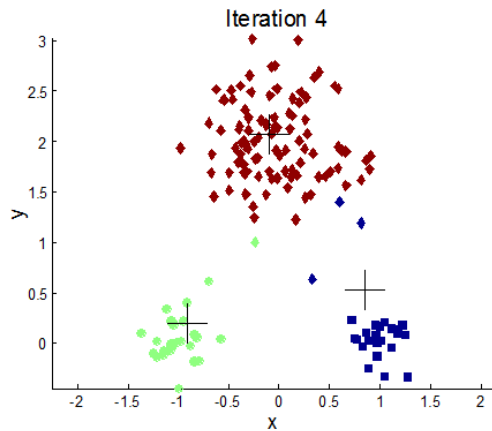
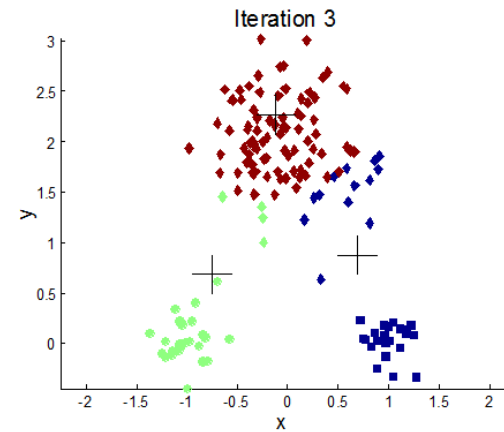
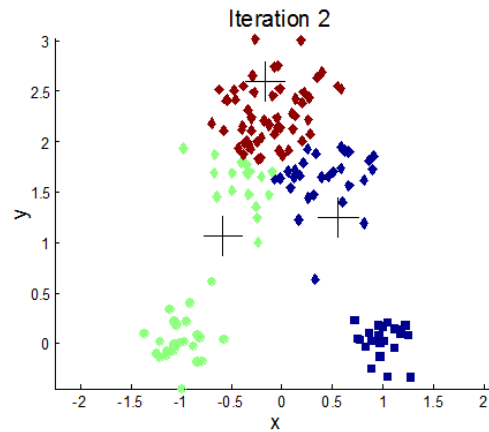
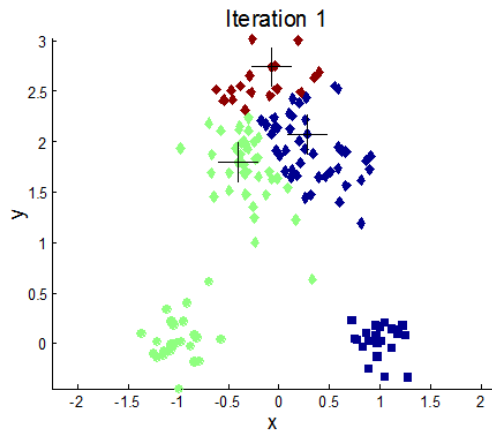
- KMC example



K-Means Clustering

- 초기 중심 설정의 영향

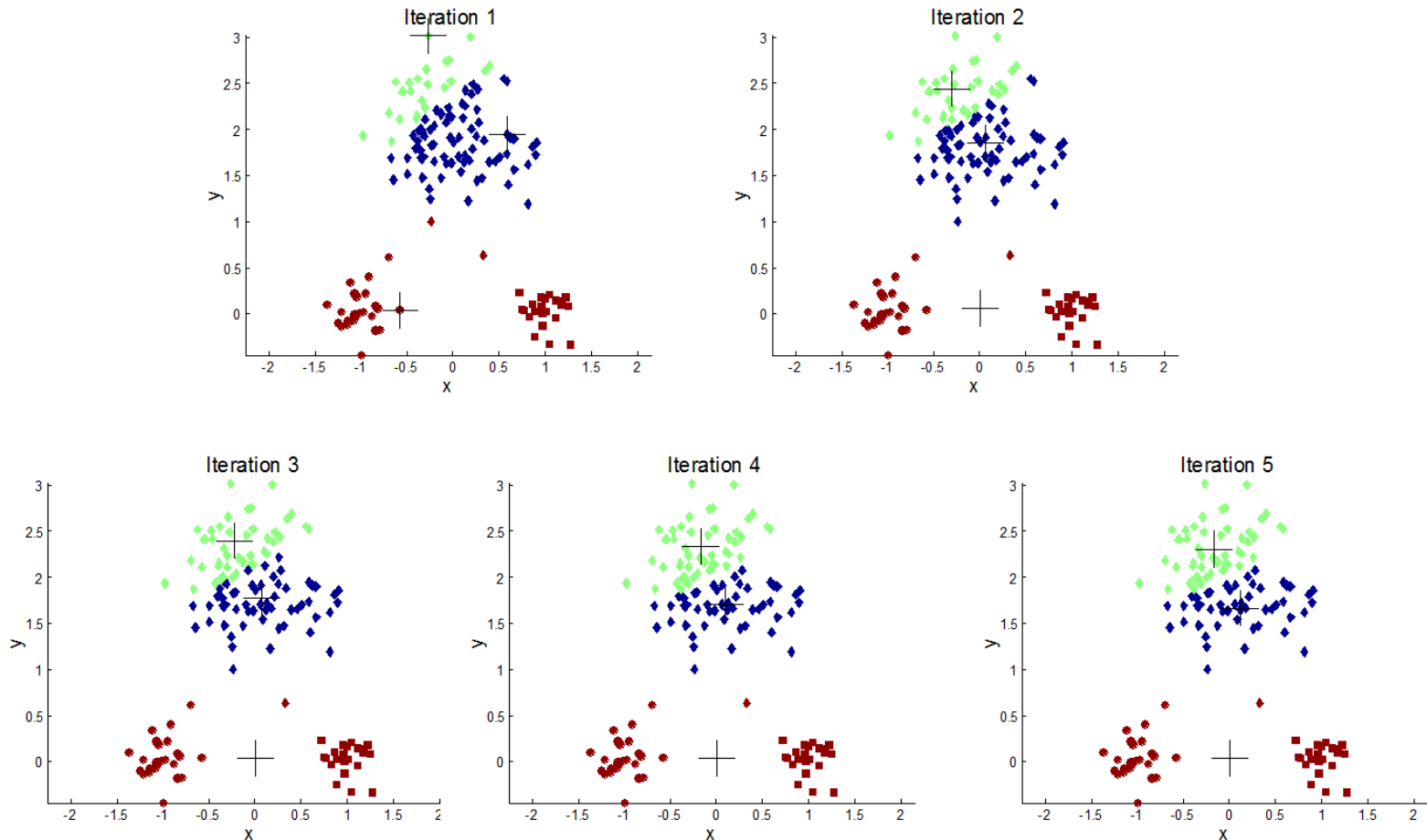
✓ 바람직한 결과



K-Means Clustering

- 초기 중심 설정의 영향

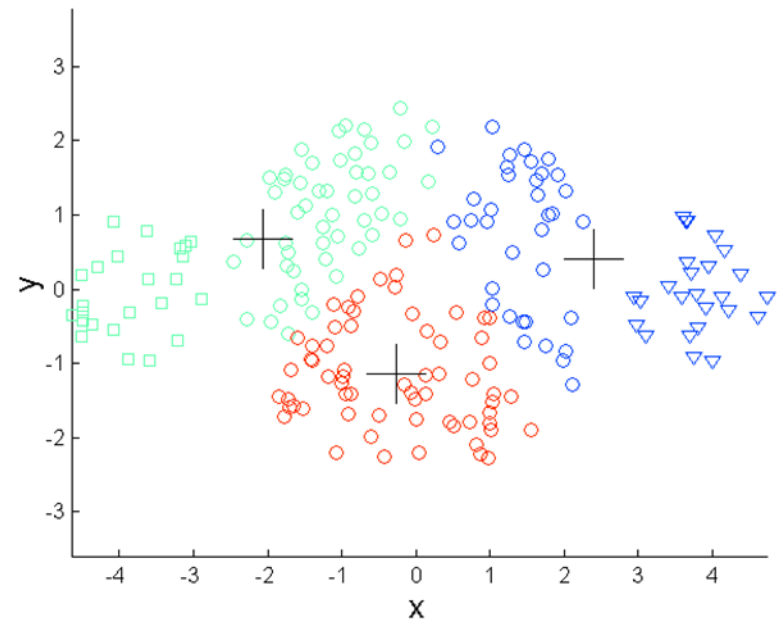
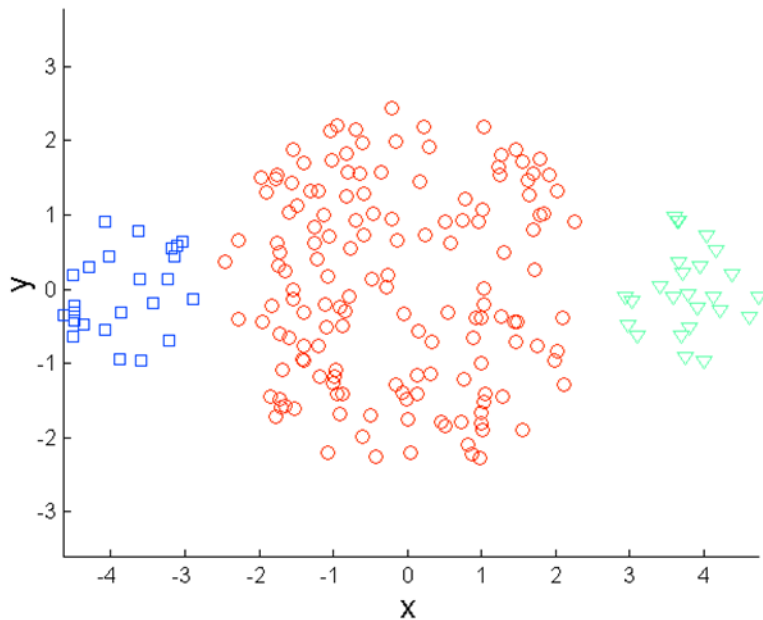
✓바람직하지 않은 결과



K-Means Clustering

- K-평균 군집화의 문제점

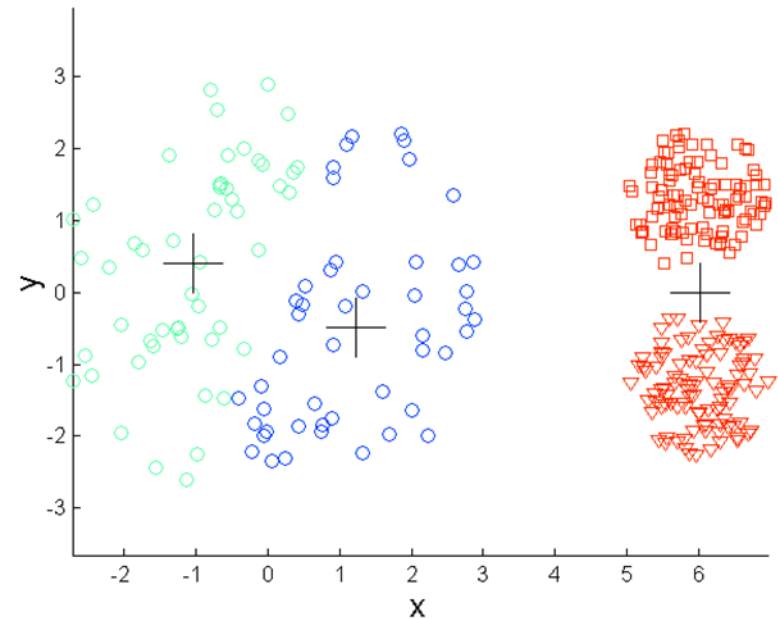
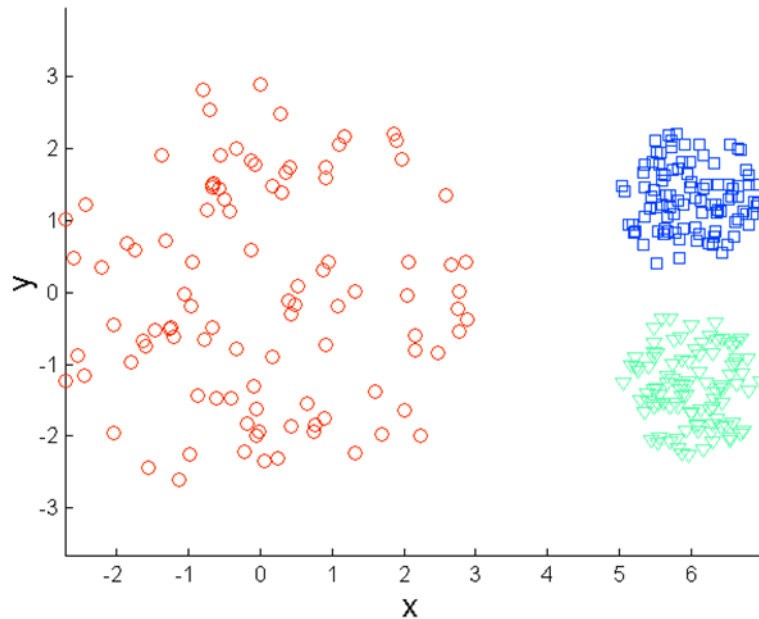
✓ 문제점 1: 서로 다른 크기의 군집을 잘 찾아내지 못함



K-Means Clustering

- K-평균 군집화의 문제점

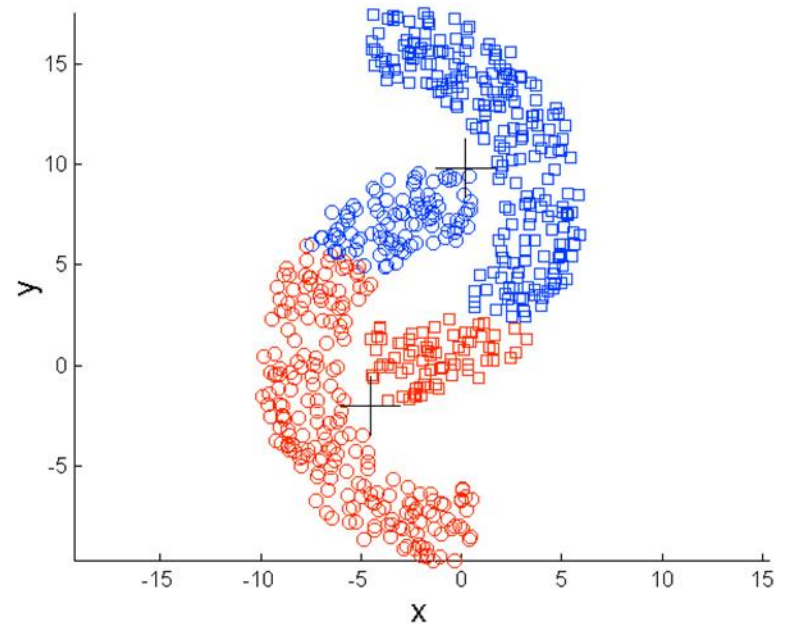
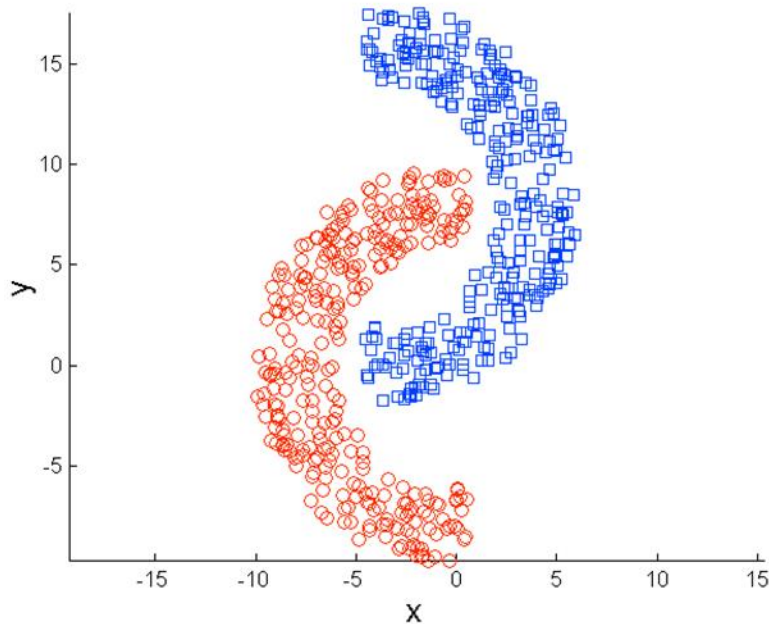
✓ 문제점 2: 서로 다른 밀도의 군집을 잘 찾아내지 못함



K-Means Clustering

- K-평균 군집화의 문제점

✓ 문제점 3: **구형이 아닌 형태**의 군집을 판별하기 어려움



AGENDA

- 01 Clustering: Overview
- 02 K-Means Clustering
- 03 Hierarchical Clustering
- 04 Density-based Clustering: DBSCAN

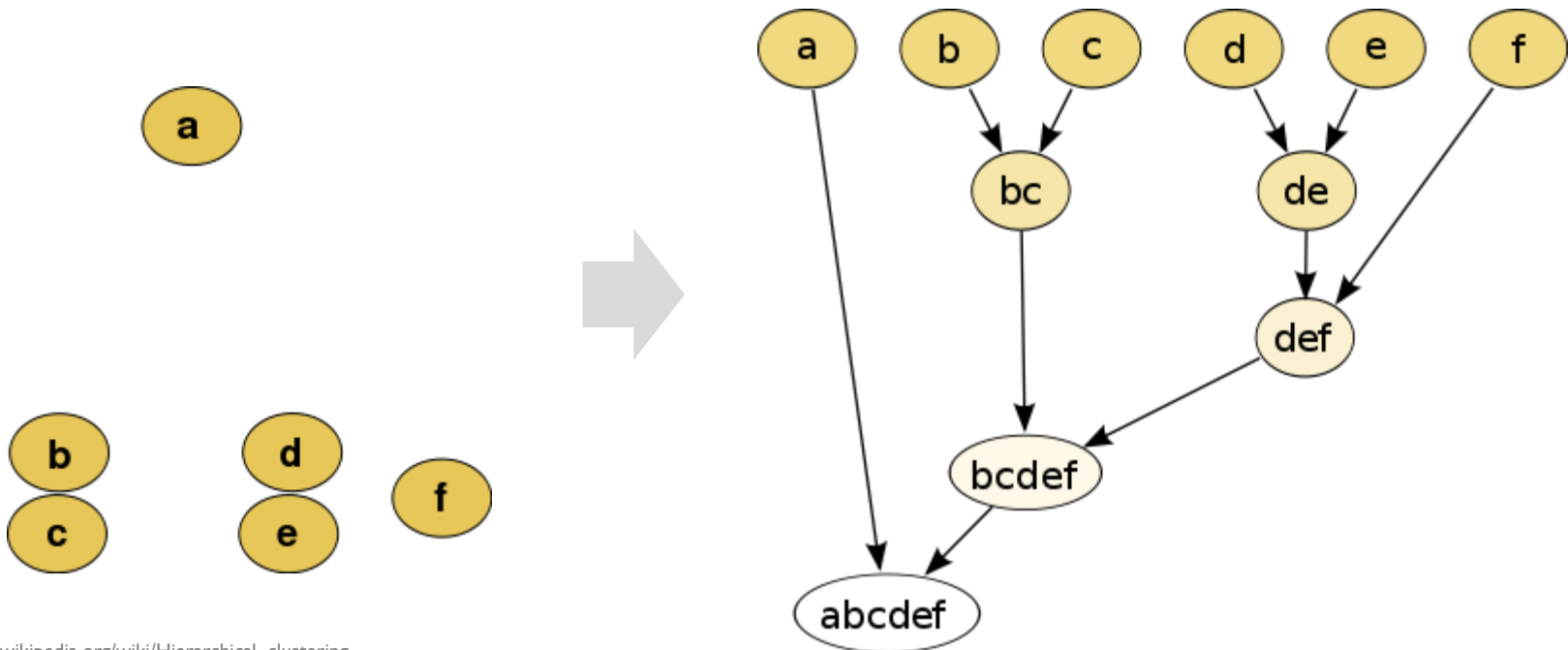
Hierarchical Clustering

- 계층적 군집화

- ✓ 계층적 트리모형을 이용하여 개별 개체들을 순차적/계층적으로 유사한 개체/군집과 통합

- ✓ 덴드로그램(Dendrogram)을 통해 시각화 가능

- Dendrogram: 개체/군집들이 결합되는 순서를 나타내는 트리형태의 구조



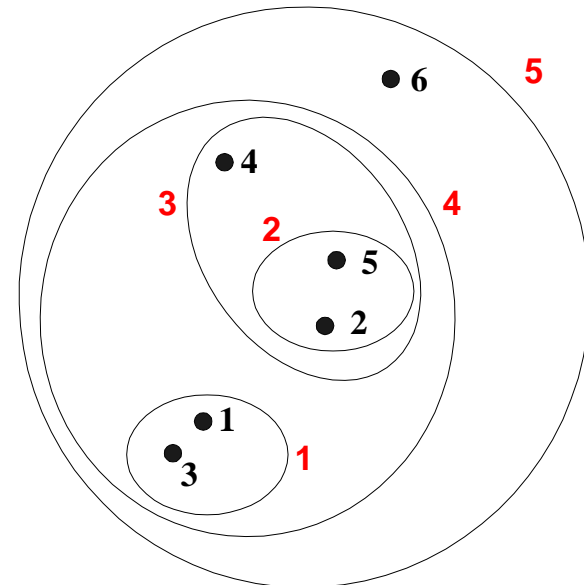
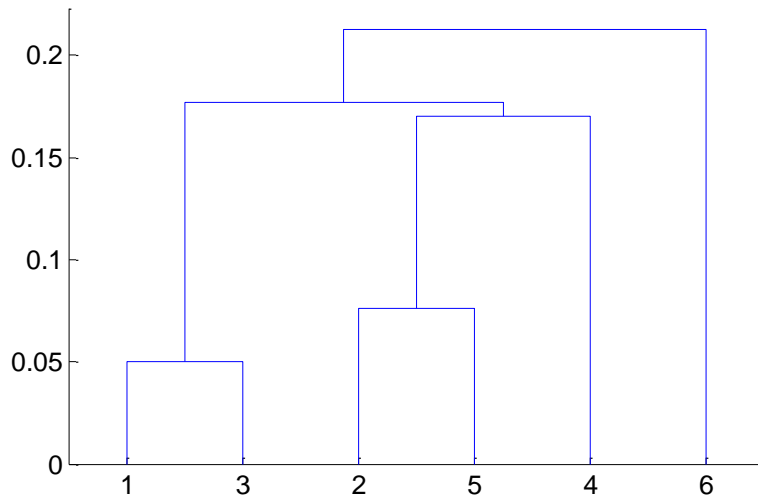
Hierarchical Clustering

- 계층적 군집화

- ✓ 계층적 트리모형을 이용하여 개별 개체들을 순차적/계층적으로 유사한 개체/군집과 통합

- ✓ 덴드로그램(Dendrogram)을 통해 시각화 가능

- Dendrogram: 개체/군집들이 결합되는 순서를 나타내는 트리형태의 구조



Hierarchical Clustering

- 계층적 군집화의 장점

- ✓ 사전에 군집의 수를 정하지 않아도 수행 가능

- Dendrogram이 생성된 후 적절한 수준에서 자르면 그에 해당하는 군집화 결과 생성

- ✓ 특정 분야(domain)에서는 이 dendrogram이 유의미한 계통체계(taxonomies)를 표현하기도 함

- 계층적 군집화의 두 가지 방식

- ✓ 상향식 군집화: Agglomerative clustering

- 초기에 모든 개체들을 개별적인 군집으로 가정

- 각 단계에서 유사한 개체/군집 결합 → 모든 개체들이 하나의 군집으로 통합되면 완료

- ✓ 하향식 군집화: Divisive clustering

- 모든 개체가 하나의 군집으로 이루어진 상태에서 출발

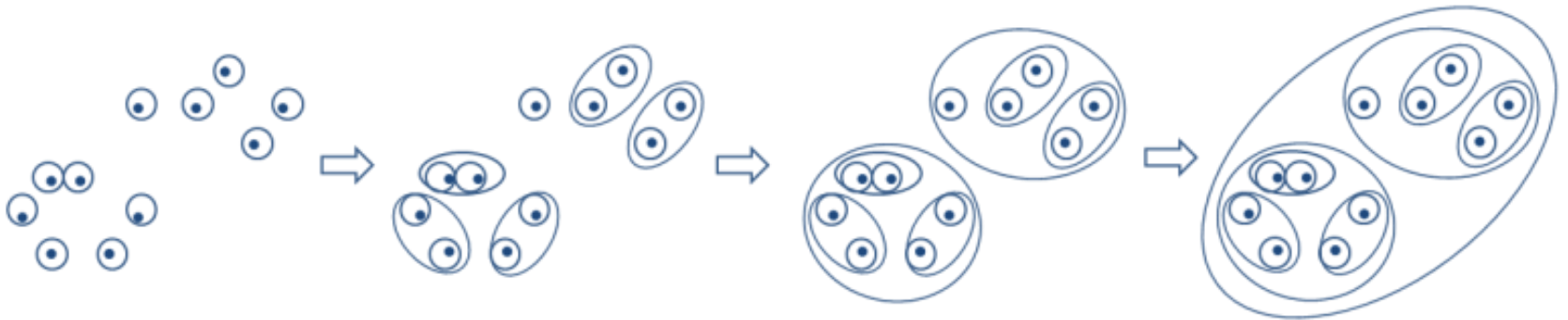
- 각 단계에서 가장 유의미하게 구분되는 지점을 판별하여 지속적으로 데이터를 분할

Hierarchical Clustering

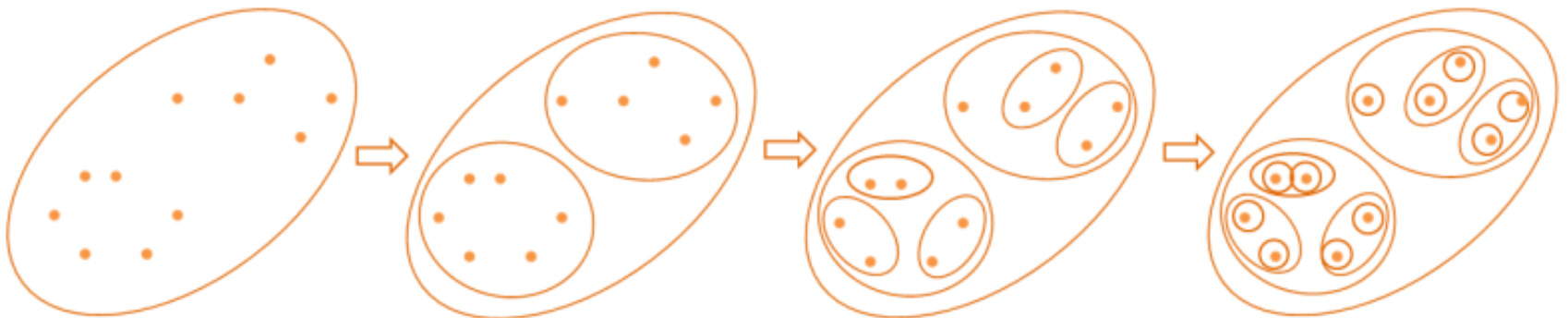
- 계층적 군집화의 두 가지 방식

- ✓ 상향식 vs. 하향식 군집화 비교

Agglomerative Hierarchical Clustering



Divisive Hierarchical Clustering

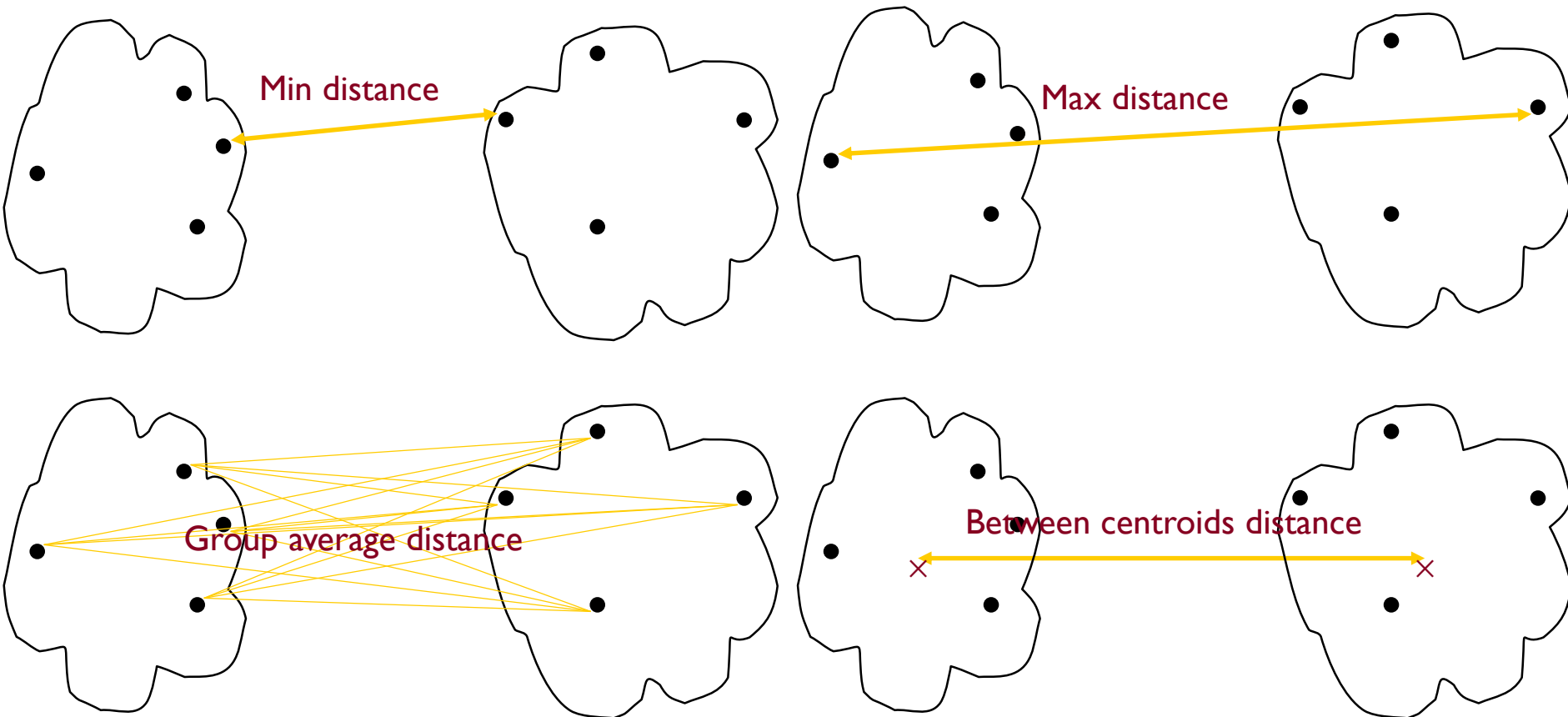


Hierarchical Clustering

- 상향식 군집화 알고리즘

- ✓ 핵심 수행 절차: 두 군집 사이의 유사도/거리 측정

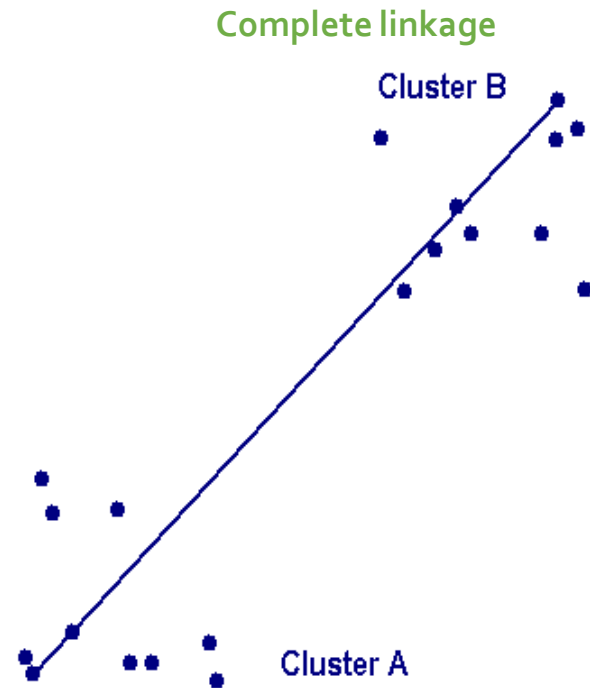
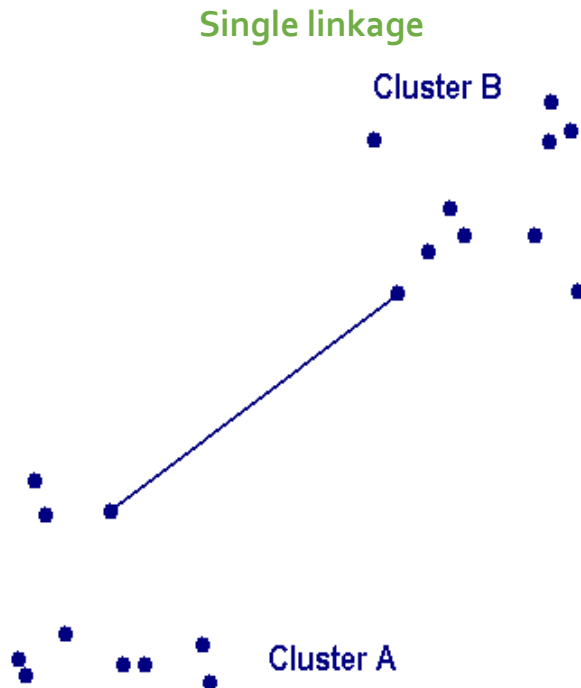
- Min, max, group average, between centroid, etc.



Hierarchical Clustering

- 상향식 군집화 알고리즘

- ✓ Single linkage (minimum distance): 각 군집에 속한 개체들 사이의 거리 중 가장 가까운 값을 군집간 거리로 정의
- ✓ Complete linkage (maximum distance): 각 군집에 속한 개체들 사이의 거리 중 가장 먼 값을 군집간 거리로 정의

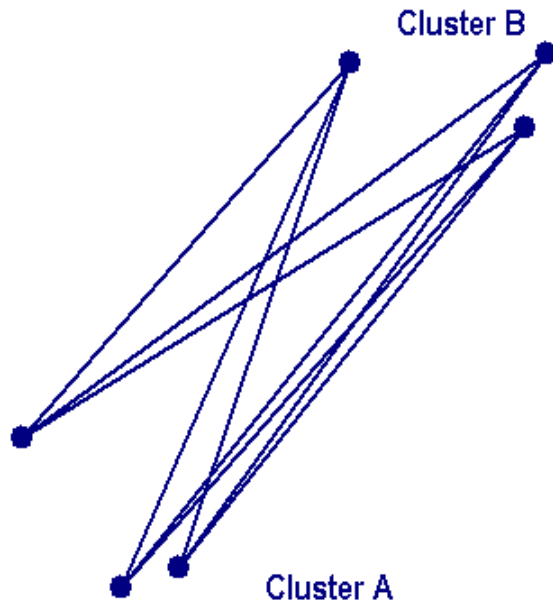


Hierarchical Clustering

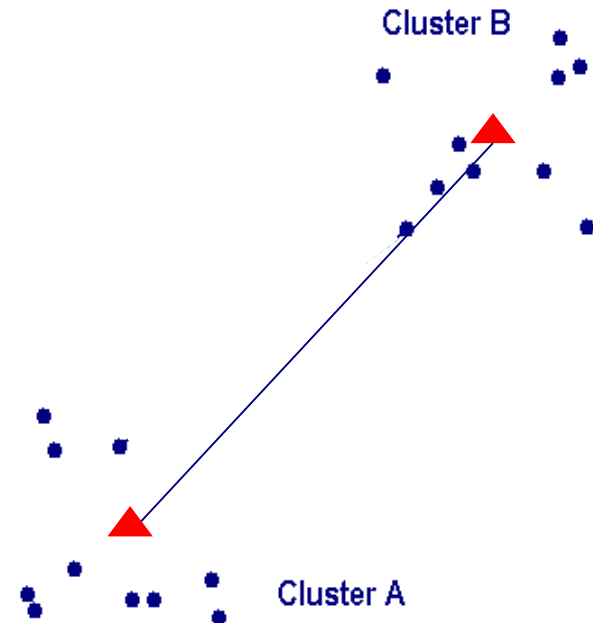
- 상향식 군집화 알고리즘

- ✓ Average linkage (mean distance): 각 군집에 속한 개체들 사이의 거리 평균값을 군집간 거리로 정의
- ✓ Centroid linkage (distance between centroids): 각 군집의 중심간 거리를 군집간 거리로 정의

Average linkage



Centroid linkage

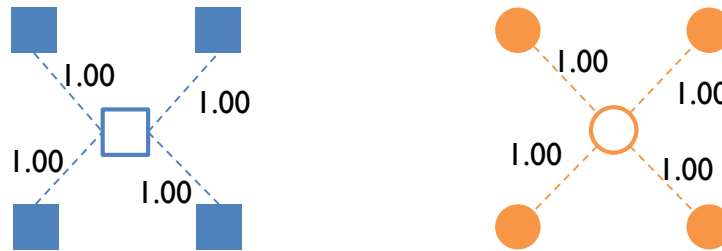


Hierarchical Clustering

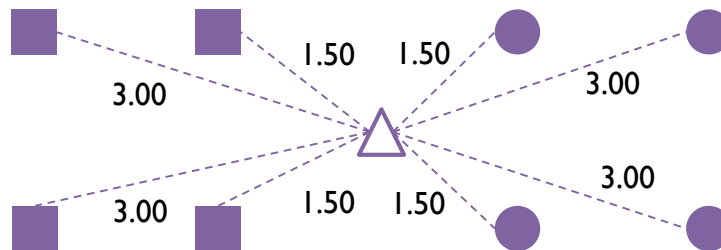
- 상향식 군집화 알고리즘: Ward's method

✓ 두 군집간의 거리를 군집이 병합된 이후의 Sum of Squared Error (SSE)와 개별 군집의 SSE의 합과의 차이로 정의

- 병합 전 SSE: $1^2 + 1^2 + 1^2 + 1^2 + 1^2 + 1^2 + 1^2 + 1^2 = 8$



- 병합 후 SSE: $4 \times 1.5^2 + 4 \times 3^2 = 45$



- Ward's distance: $45 - 8 = 37$

Hierarchical Clustering

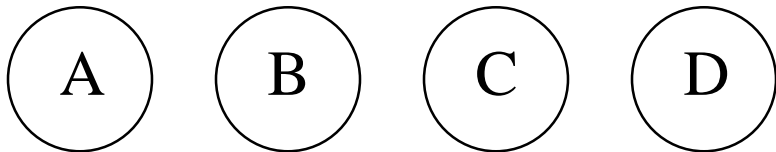
- 상향식 군집화 알고리즘

- ✓ 1단계: 모든 개체를 개별 군집으로 정의하고 군집간 거리 행렬 계산
- ✓ 2단계: 다음 절차를 반복
 - 2-1단계: 가장 가까운 두 개의 군집을 하나의 군집으로 통합
 - 2-2단계: 군집간 거리 행렬 업데이트
- ✓ 종료 조건: 모든 개체가 하나의 군집으로 통합되면 종료

Hierarchical Clustering

- 계층적 군집화 절차 예시

Initial Data Items



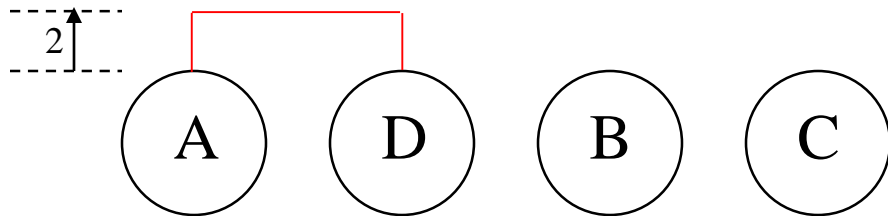
Distance Matrix

Dist	A	B	C	D
A		20	7	2
B			10	25
C				3
D				

Hierarchical Clustering

- 계층적 군집화 절차 예시

Current Clusters



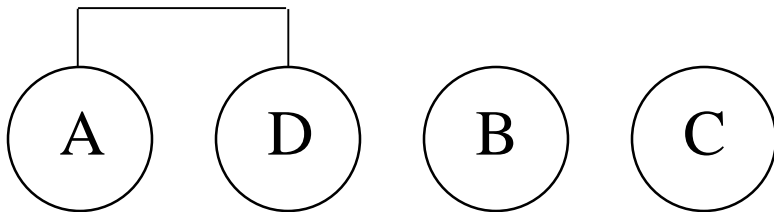
Distance Matrix

Dist	A	B	C	D
A		20	7	2
B			10	25
C				3
D				

Hierarchical Clustering

- 계층적 군집화 절차 예시

Current Clusters



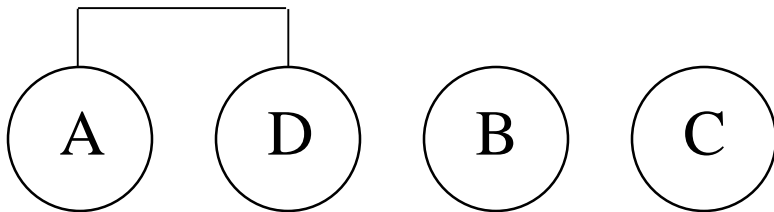
Distance Matrix

Dist	AD	B	C	
AD		20	3	
B			10	
C				

Hierarchical Clustering

- 계층적 군집화 절차 예시

Current Clusters



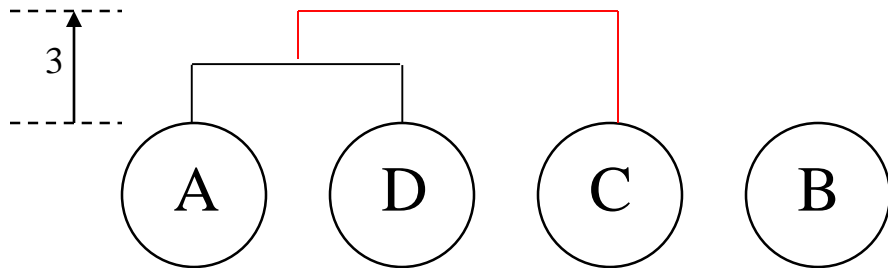
Distance Matrix

Dist	AD	B	C	
AD		20	3	
B			10	
C				

Hierarchical Clustering

- 계층적 군집화 절차 예시

Current Clusters



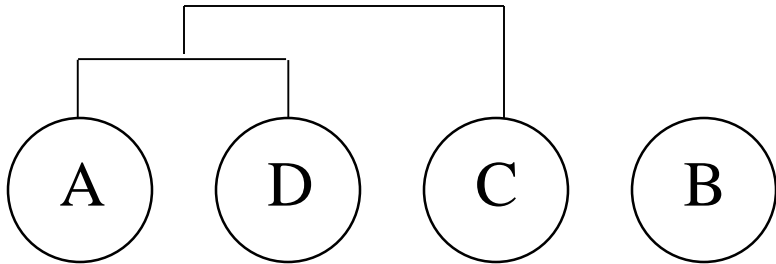
Distance Matrix

Dist	AD	B	C	
AD		20	3	
B			10	
C				

Hierarchical Clustering

- 계층적 군집화 절차 예시

Current Clusters



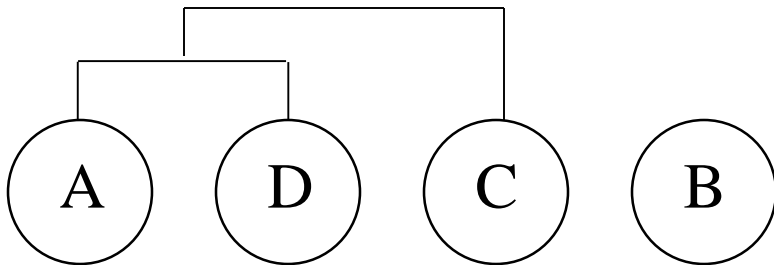
Distance Matrix

Dist	AD C	B		
AD C		10		
B				

Hierarchical Clustering

- 계층적 군집화 절차 예시

Current Clusters



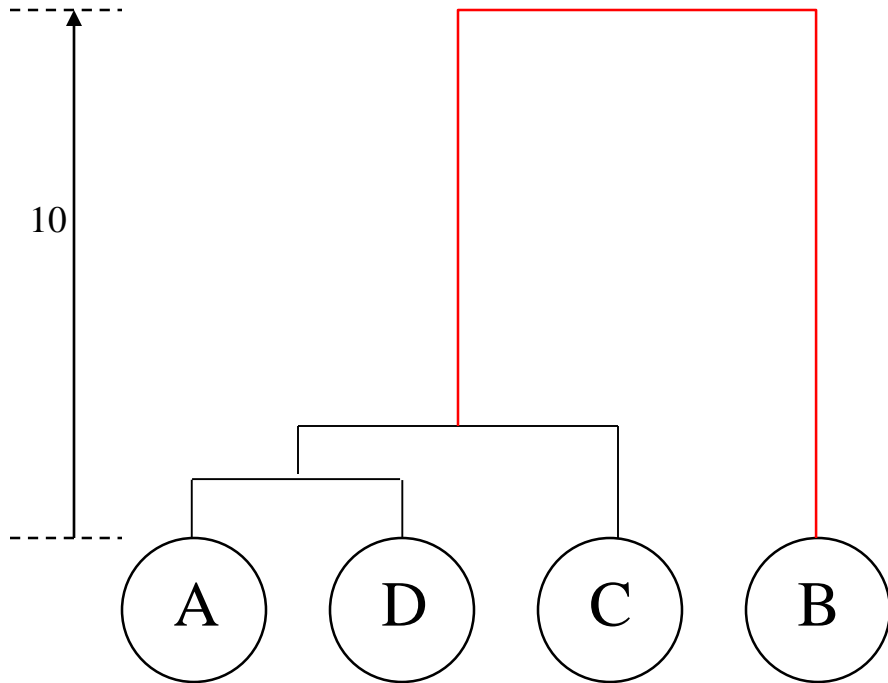
Distance Matrix

Dist	AD C	B		
AD C		10		
B				

Hierarchical Clustering

- 계층적 군집화 절차 예시

Current Clusters



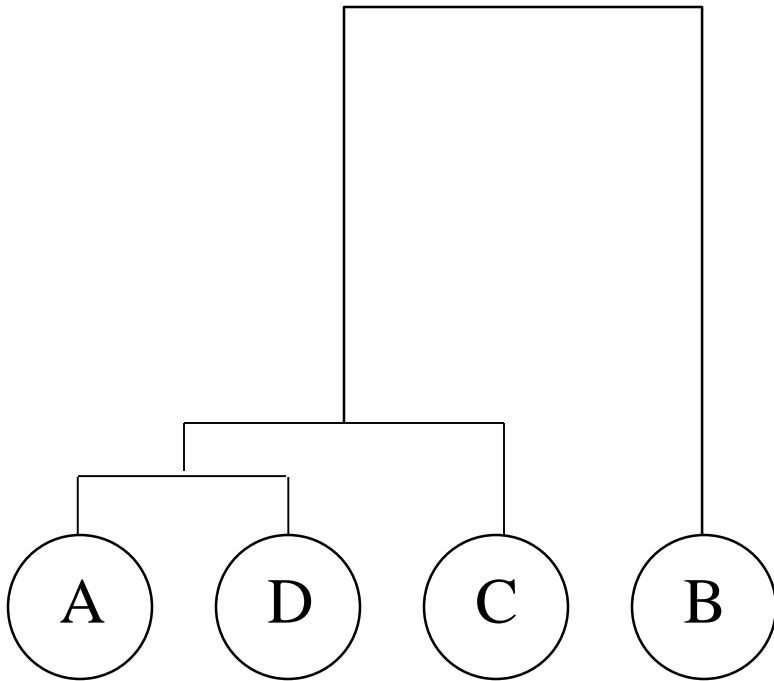
Distance Matrix

Dist	AD C	B		
AD C		10		
B				

Hierarchical Clustering

- 계층적 군집화 절차 예시

Final Result

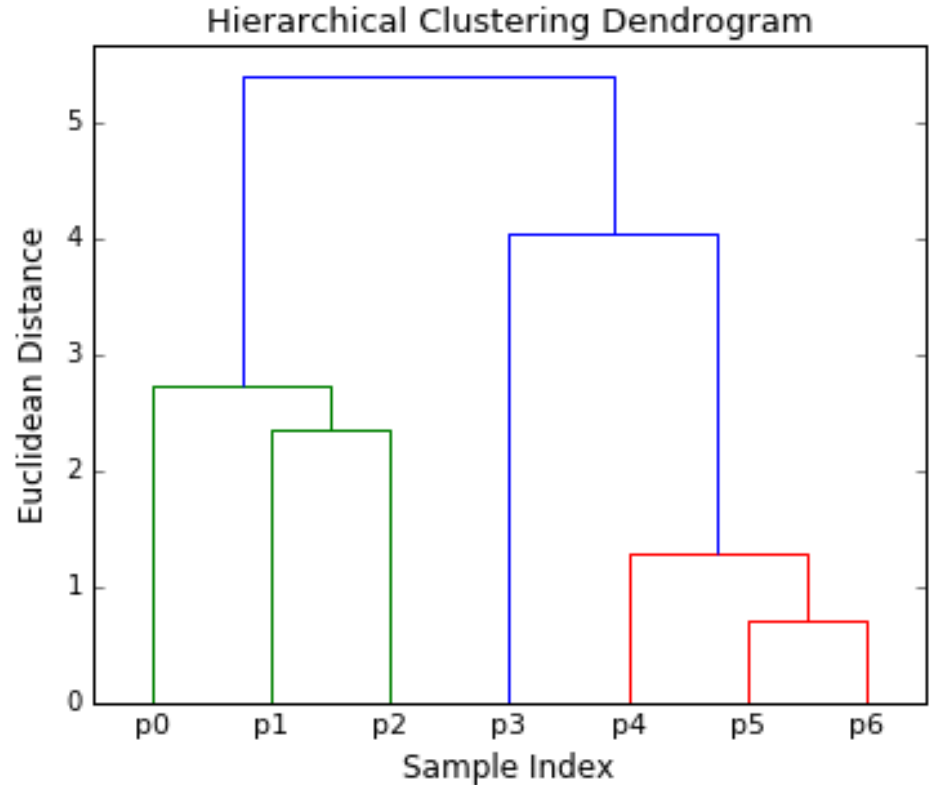
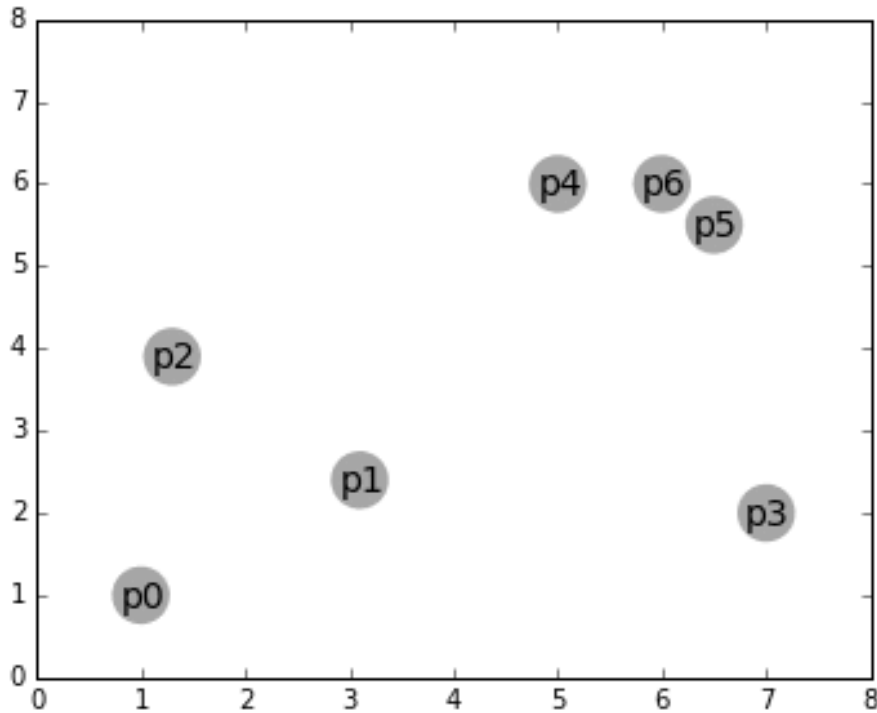


Distance Matrix

Dist	AD CB			
AD CB				

Hierarchical Clustering

- 계층적 군집화 절차 예시



<https://towardsdatascience.com/the-5-clustering-algorithms-data-scientists-need-to-know-a36d136ef68>

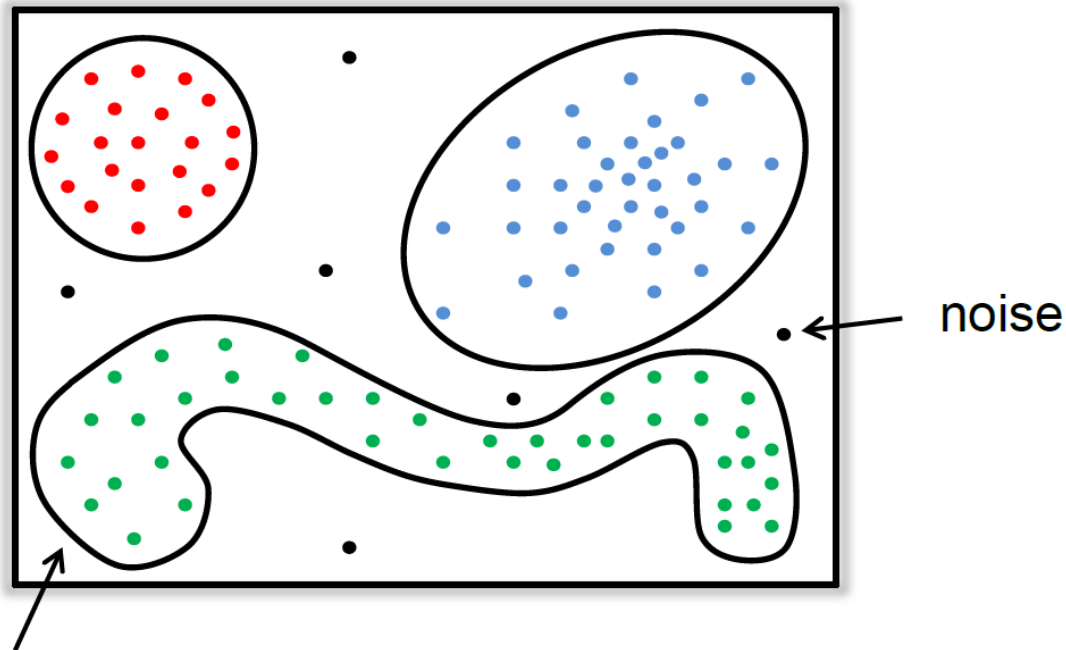
AGENDA

- 01 Clustering: Overview
- 02 K-Means Clustering
- 03 Hierarchical Clustering
- 04 **Density-based Clustering: DBSCAN**

Density-based Clustering

Ester et al. (1996)

- Density-based clustering
 - ✓ Conduct a clustering by considering the density of data points
 - Can find an arbitrary shape of cluster
 - Can remove noise from clustering result



arbitrarily shaped clusters

Density-based Clustering

- DBSCAN

- ✓ Most popular density-based clustering algorithm

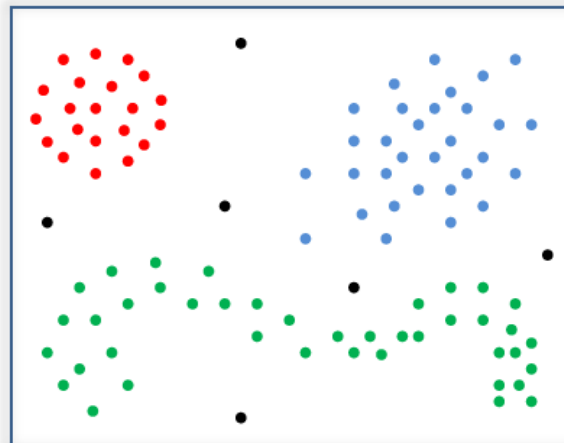
- Idea

- ✓ Clusters are the collections of data points with high density

- ✓ Density around a noise point is very low

- Purpose

- ✓ Quantify the features of clusters and noise points to find a set of valid clusters



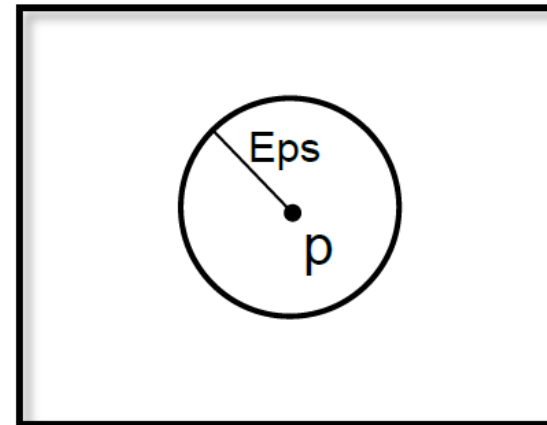
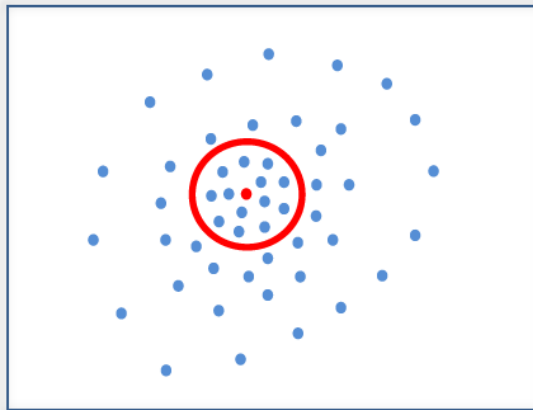
Density-based Clustering

- DBSCAN

- ✓ Definition 1: ϵ -neighborhood of a point

- The ϵ -neighborhood of a point, denoted by $N_\epsilon(p)$, is defined by

$$N_\epsilon(p) = \{q \in D \mid \text{dist}(p, q) \leq \epsilon\}$$



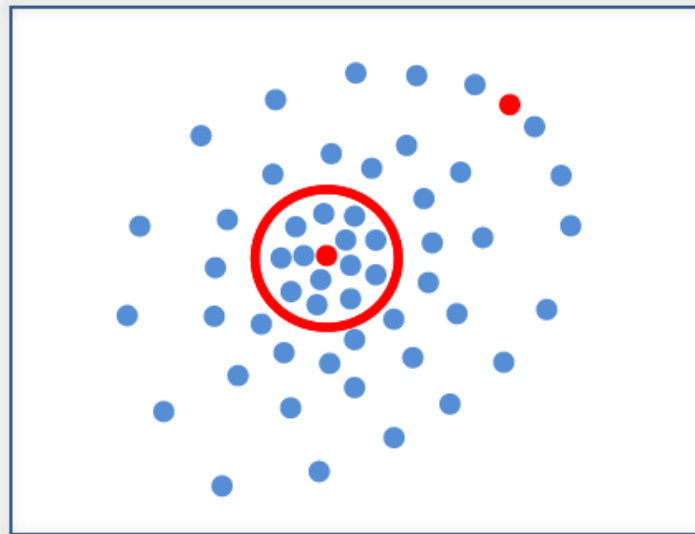
- ✓ Naïve Approach: require for each point in a cluster that there are at least a minimum number (MinPts) of points in an ϵ -neighborhood of that point

Density-based Clustering

- DBSCAN

- ✓ Problem of Naïve Approach

- There are two kinds of points in a cluster
 - Points inside of the cluster (core points)
 - Points on the border of the cluster (border points)
 - An ϵ -neighborhood of a border point contains significantly less points than an ϵ -neighborhood of a core point

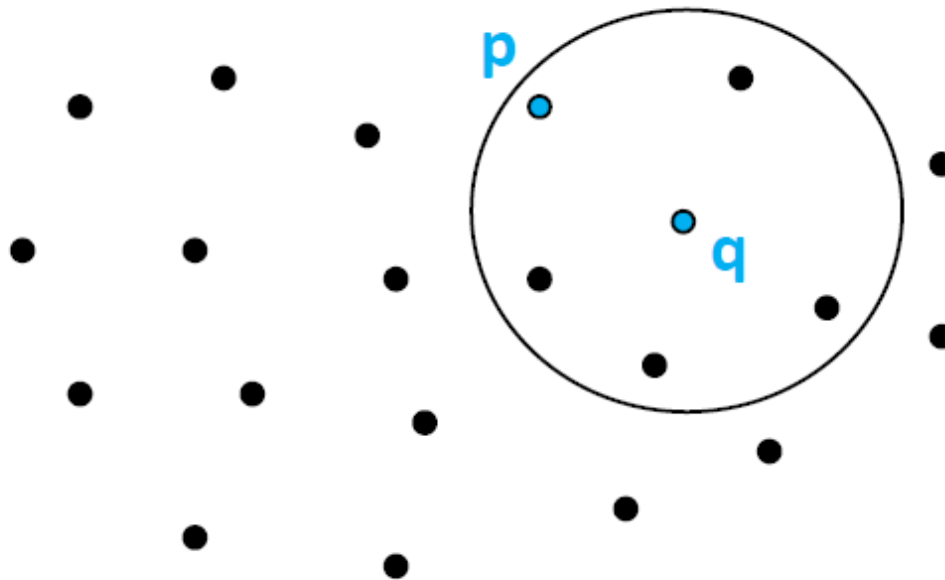


Density-based Clustering

- DBSCAN

- ✓ Better idea

- For every point p in a cluster C , there is a point q in C so that p is inside of the ε -neighborhood of q (Border points are connected to core points)
 - $N_\varepsilon(q)$ contains at least MinPts points (Core points = high density)



Density-based Clustering

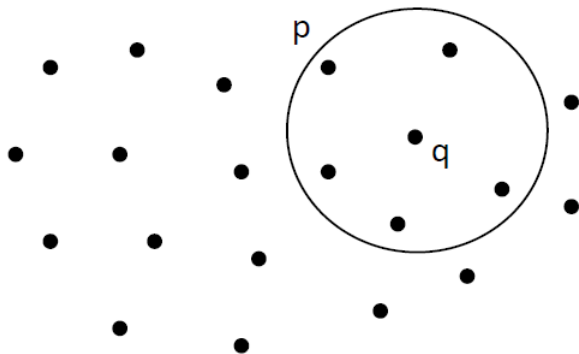
- DBSCAN

- ✓ Definition 2: directly density-reachable

- A point p is directly density-reachable from a point q with regard to the parameters ϵ and $MinPts$, if

1) $p \in N_\epsilon(q)$ (*reachability*)

2) $|N_\epsilon(q)| \geq MinPts$ (*core point condition*)



$MinPts = 5$

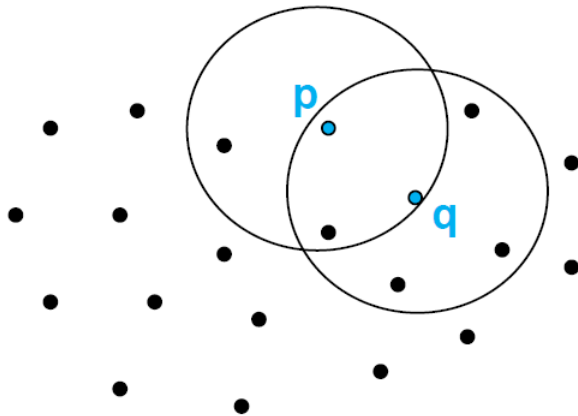
$|N_{Eps}(q)| = 6 \geq 5 = MinPts$ (core point condition)

Density-based Clustering

- DBSCAN

- ✓ Property

- Directly density-reachable is symmetric for **pairs of core points**
 - It is not symmetric if **one core point** and **one border point** are involved



Parameter: MinPts = 5

p directly density reachable from q

$$p \in N_{Eps}(q)$$

$$|N_{Eps}(q)| = 6 \geq 5 = \text{MinPts} \quad (\text{core point condition})$$

q not directly density reachable from p

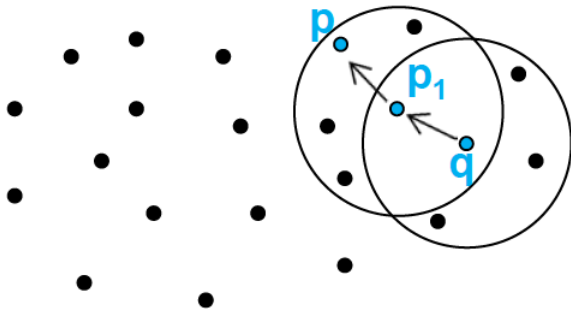
$$|N_{Eps}(p)| = 4 < 5 = \text{MinPts} \quad (\text{core point condition})$$

Density-based Clustering

- DBSCAN

- ✓ Definition 3: density-reachable

- A point p is density-reachable from a point q with regard to the parameters ϵ and MinPts , if there is a chain of points p_1, p_2, \dots, p_s with $p_1 = q$ and $p_s = p$ such that p_{i+1} is directly density-reachable from p_i for all $1 < i < s-1$



$\text{MinPts} = 5$

$|N_{\text{Eps}}(q)| = 5 = \text{MinPts}$ (core point condition)

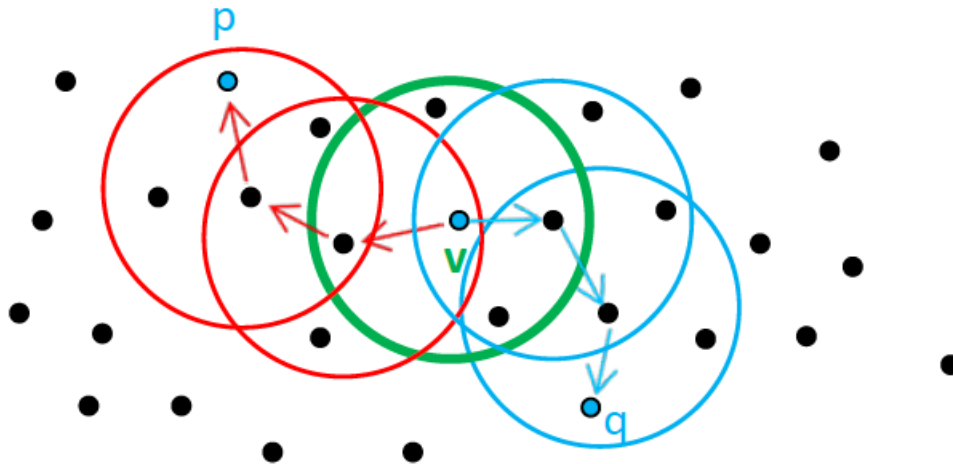
$|N_{\text{Eps}}(p_1)| = 6 \geq 5 = \text{MinPts}$ (core point condition)

Density-based Clustering

- DBSCAN

- ✓ Definition 4: density-connected

- A point p is density-connected to a point q with regard to the parameters ϵ and MinPts, if there is a point v such that both p and q are density-reachable from v



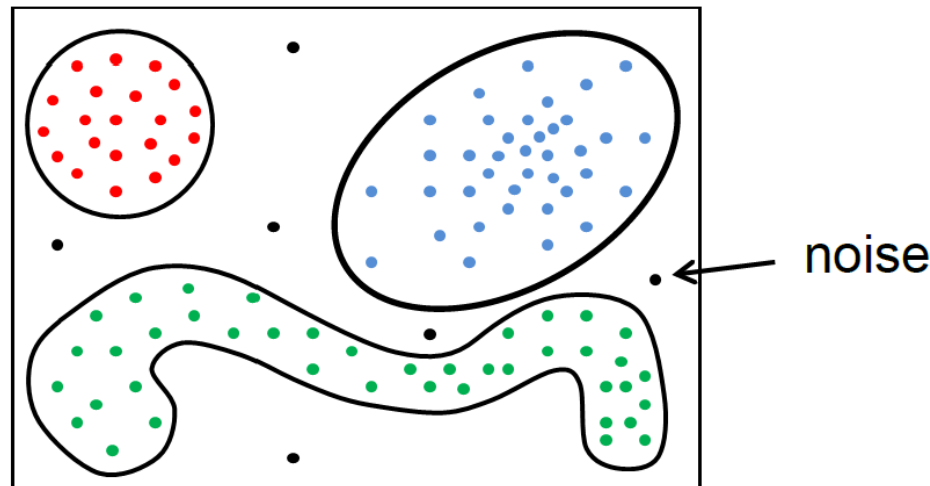
MinPts = 5

Density-based Clustering

- DBSCAN

- ✓ Definition 5: Cluster

- A cluster with regard to the parameters ϵ and MinPts is a non-empty subset C of the database D with
 - (1) For all $p, q \in D$: If $p \in C$ and q is density-reachable from p with regard to the parameters ϵ and MinPts, then $q \in C$ (**Maximality**)
 - (2) For all $p, q \in C$: The point p is density-connected to q with regard to the parameters ϵ and MinPts (**Connectivity**)

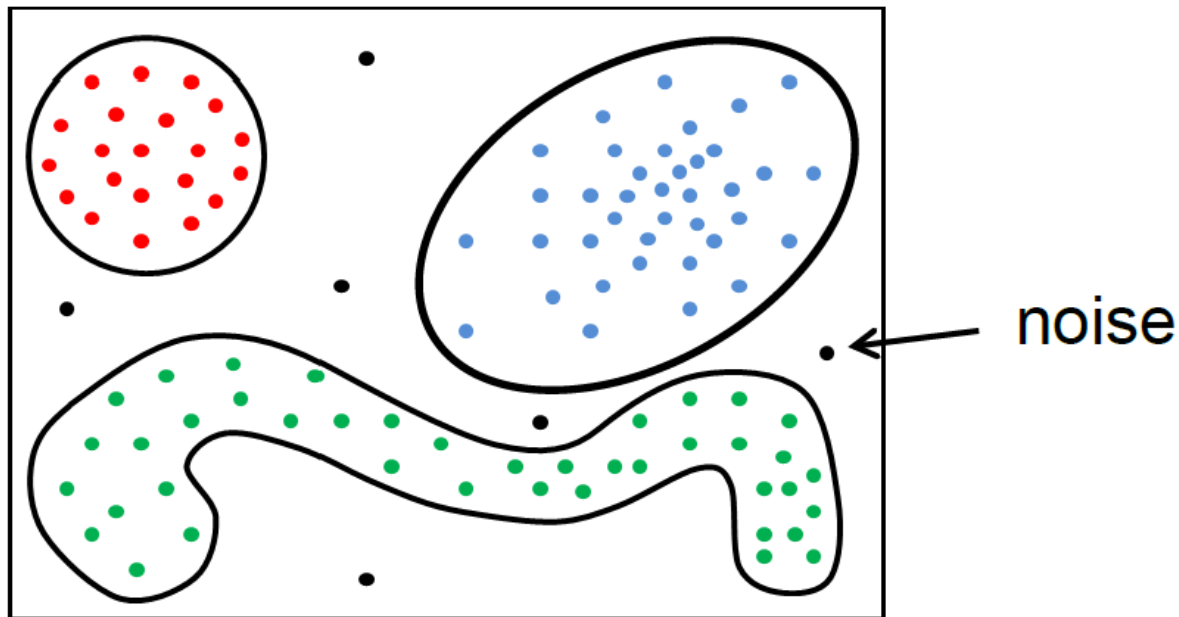


Density-based Clustering

- DBSCAN

- ✓ Definition 6: Noise

- Let C_1, \dots, C_k be the clusters of the database D with regard to the parameters ε and MinPts
 - The set of points in the database D not belonging to any cluster C_1, \dots, C_k is called noise



Density-based Clustering

- DBSCAN: Algorithm

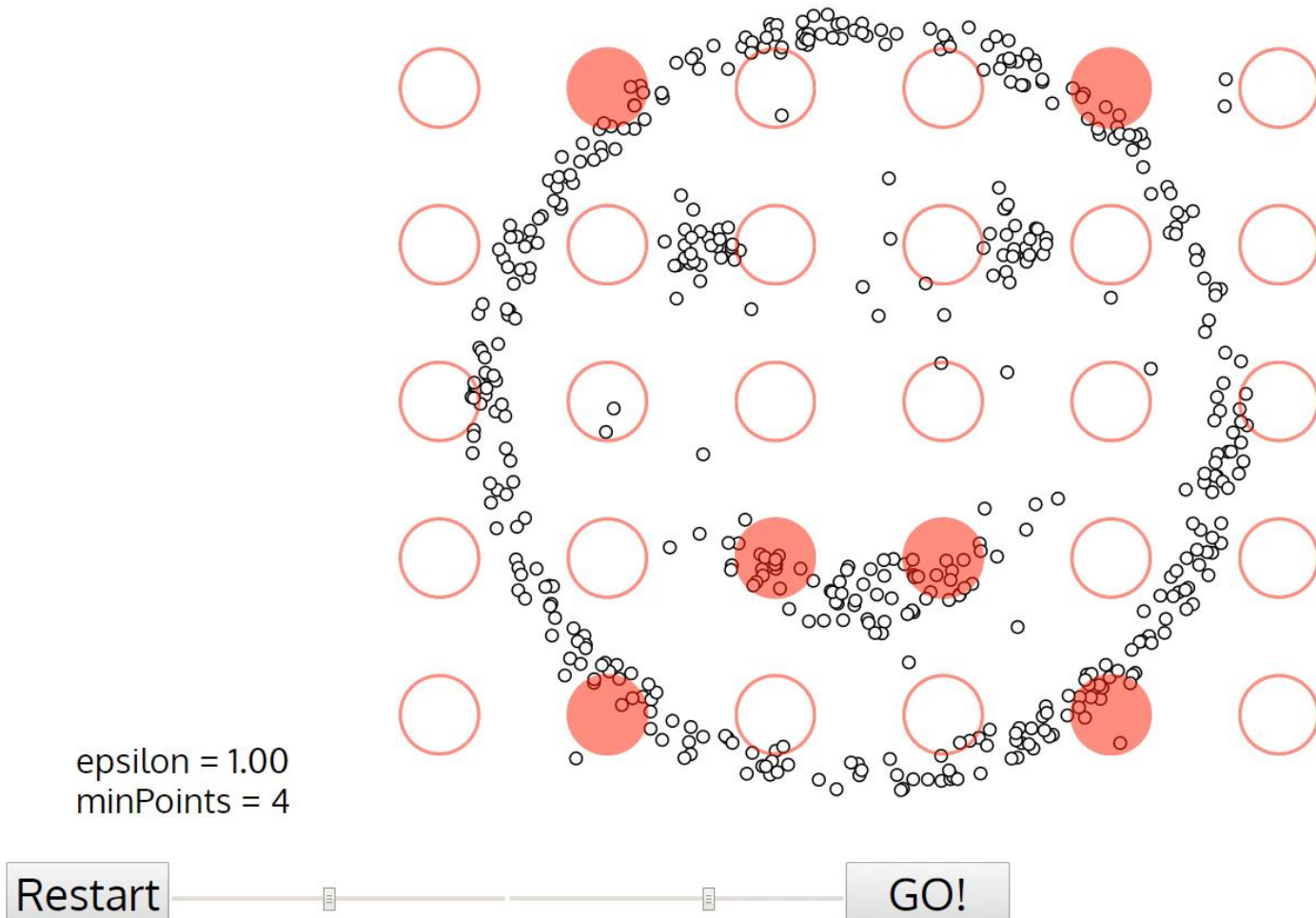
- ✓ Input: N objects to be clustered and global parameter, ϵ and MinPts
- ✓ Output: Cluster of objects

- Algorithm

- ✓ Arbitrary select a point p
- ✓ Retrieve all points density-reachable from p w.r.t. ϵ and MinPts
- ✓ If p is a core points, 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

Density-based Clustering

- DBSCAN example

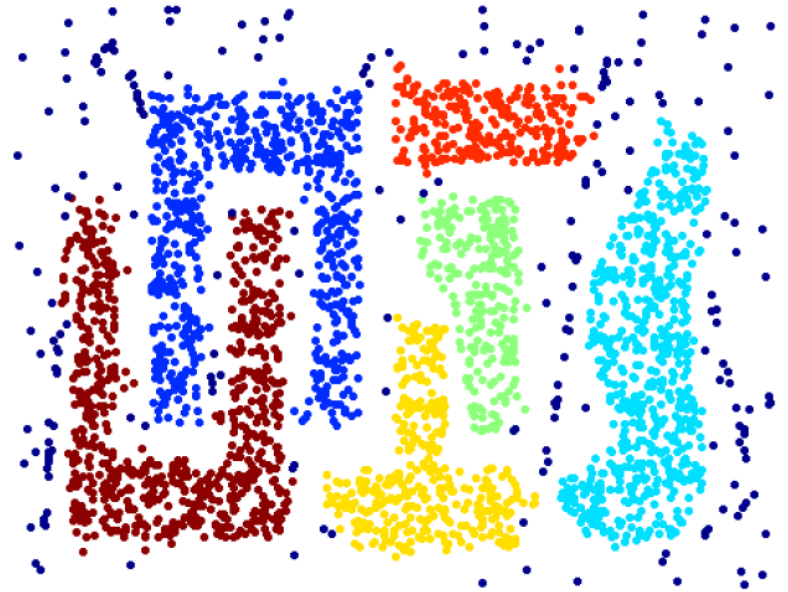


Density-based Clustering

- DBSCAN example



Original Points



Clusters

