

✓ 군집분석이란?

- 데이터 간의 유사도를 정의하고, 그 유사도에 가까운 것부터 순서대로 합쳐 가는 방법으로 그룹(군집)을 형성한 후 각 그룹의 성격을 파악하거나 그룹 간의 비교분석을 통해서 데이터 전체의 구조에 대한 이해를 돕고자 하는 탐색적인 분석 방법

군집 분석의 목적

- 데이터 셋 전체를 대상으로 서로 유사한 개체 들을 몇 개의 군집으로 세분화하여 대상 집단을 정확하게 이해하고, 효율적으로 활용하기 위함.
 - 군집 분석으로 그룹화된 군집은 변수의 특성이 그룹 내적으로는 동일하고, 외적으로는 이질적인 특성을 가짐
 - 군집 분석의 용도는 고객의 충성도에 따라서 몇 개의 그룹으로 분류하고, 그룹별로 맞춤형 마케팅 및 프로모션 전략을 수립하는 데 활용됨

계층적 군집분석

- 데이터를 유사도에 따라 계층적으로 묶어주는 군집분석 방법
- 예를 들어, 고객 데이터를 구매 패턴과 구매 금액에 따라 묶어주는 것이나, 상품 데이터를 재료, 맛, 가격 등에 따라 묶어주는 것이 대표적인 예임
- 이렇게 묶인 군집들은 덴드로그램 형태로 시각화되어 나타남

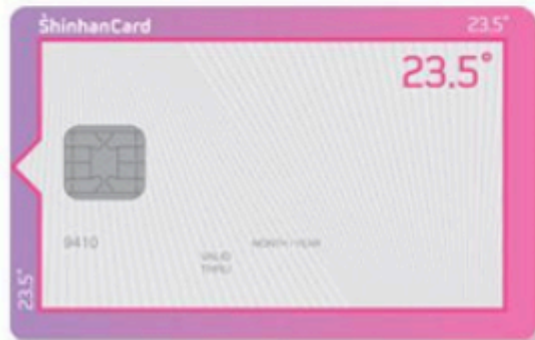
비계층적 군집분석(Non-Hierarchical Clustering)

- 계층 구조를 형성하지 않고 데이터 포인트들을 클러스터로 그룹화하는 방법으로, 주로 K-means와 같은 알고리즘이 사용됨
 - K-means의 경우, 미리 정해진 클러스터 개수(K)를 기준으로 중심점을 설정하고, 각 데이터 포인트가 가장 가까운 중심점에 속하게 클러스터링하는 방식으로 진행됨
- 비계층적 군집분석 사례:

- 마케팅 - 고객 데이터를 분석하여 고객 세분화(customer segmentation)를 수행하고, 각 세분화된 고객 그룹에 맞춘 마케팅 전략을 수립할 수 있음
- 이미지 분석 - 이미지 픽셀을 클러스터링하여 각 클러스터의 대표 색상을 추출하여 이미지 색상 분석이나 색상 분류를 할 수 있음
- 유전자 데이터 분석 - 유전자 데이터에서 유전자 클러스터를 식별하고, 유전자 간의 관련성을 분석할 수 있음
- 문서 데이터 분석 - 문서 데이터를 분석하여 유사한 내용을 가진 문서들을 클러스터링하여 문서 분류를 할 수 있음
- 고객 행동 분석 - 웹사이트 방문 고객의 행동 패턴을 분석하여 고객 세분화 및 맞춤형 서비스 제공에 활용할 수 있음

✓ 군집분석 사례 - 신한카드(고객들을 군집 분석해서 고객마다 원하는 카드혜택이 다르게 하고, 고객 특성을 새롭게 정의함으로써 고객 특성별로 특화된 카드 상품을 만드는 데 활용)

- 신한카드는 2014년부터 빅데이터를 활용하여 고객의 소비 특성을 파악하여 고객 분류를 하고, 각 고객 그룹에 맞는 카드 혜택과 유형을 정했음
- 그래서 나온 것이 바로 Code 9(코드나인)이라고 하는 것임
- 고객정보와 카드 결제내역을 군집분석하여 남녀 고객을 위와 같이 총 18가지로 분류하고 각 분류별로 Rookie, It-Girl 같은 명칭을 지었음
- 각 고객군집의 특성을 파악해서 카드 이름도 지은 것 같음



10강. 분할적 군집분석

- K-Means와 DBSCAN
- 퍼지군집화 / EM 알고리즘 / SOM
- 엘보우 기법

■ 비지도학습의 종류



■ 분할적 군집분석

계층적 관계가 없는 다수의 군집들을 만드는 방법

중심점 기반

· K-Means 군집화

밀도 기반

· DBSCAN 군집화

확률 기반

· 퍼지군집화

분포 기반

· EM알고리즘

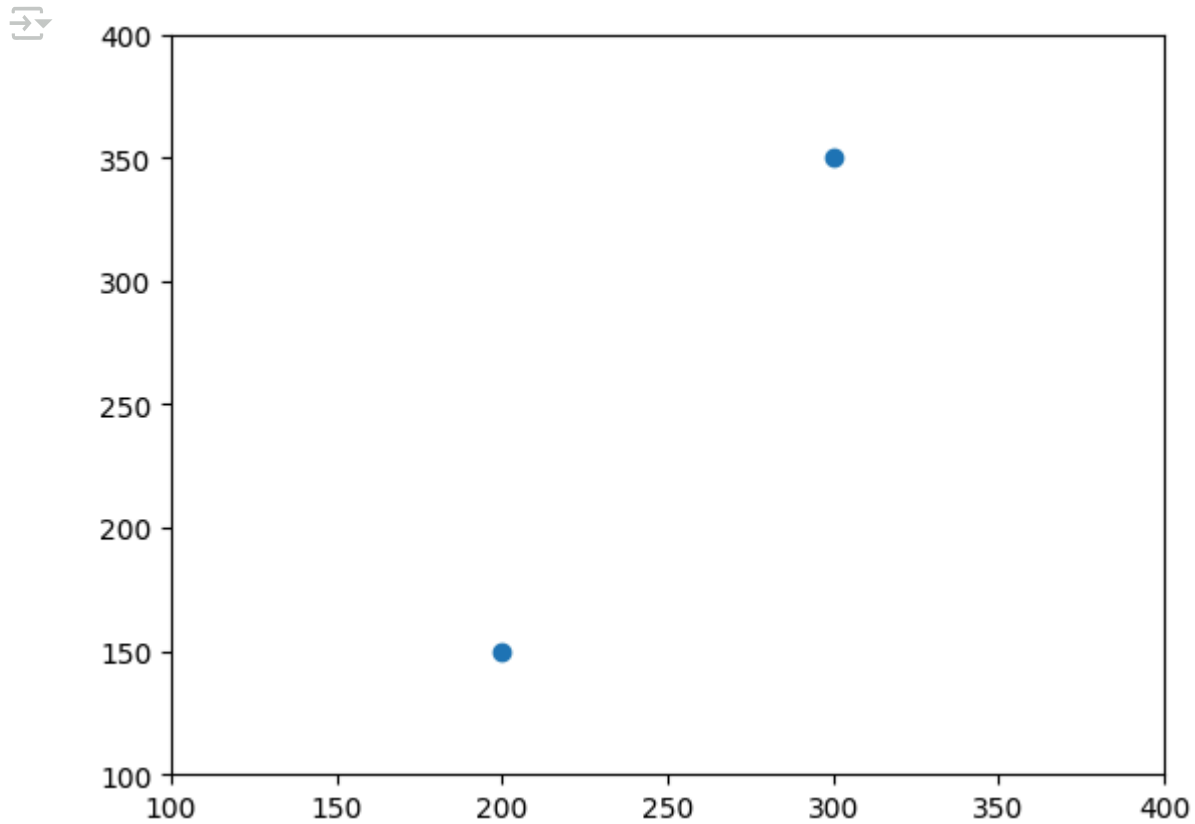
그래프 기반

· 자기조직화지도(SOM)

■ K-Means 군집화

각 데이터와 중심점의 거리를 측정 후 가장 가까운 그룹에 할당하여 K개의 군집으로 묶는 방법

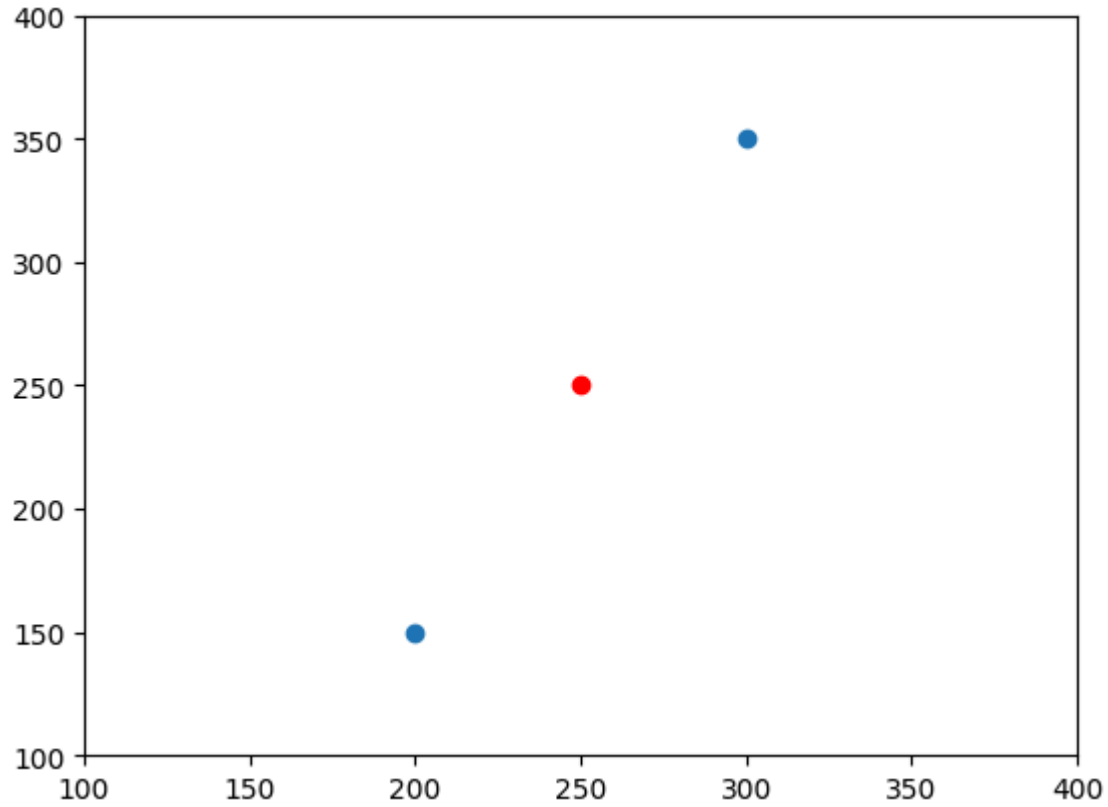
```
1 import matplotlib.pyplot as plt
2
3 x1, y1 = 300, 350
4 x2, y2 = 200, 150
5
6 plt.xlim(100, 400)
7 plt.ylim(100, 400)
8 plt.scatter([x1, x2], [y1, y2])
9 plt.show()
```



```
1 center_coordinate_x = (x1 + x2)/2
2 center_coordinate_y = (y1 + y2)/2
3 print(center_coordinate_x, center_coordinate_y)
```

```
4 plt.xlim(100, 400)
5 plt.ylim(100, 400)
6 plt.scatter(center_coordinate_x, center_coordinate_y, color = 'red')
7 plt.scatter([x1, x2], [y1, y2])
8 plt.show()
```

↔ 250.0 250.0



K-Means 군집화

각 데이터와 중심점의 거리를 측정 후 가장 가까운 그룹에 할당하여 K개의 군집으로 묶는 방법



K-Means 군집화

각 데이터와 중심점의 거리를 측정 후 가장 가까운 그룹에 할당하여 K개의 군집으로 묶는 방법



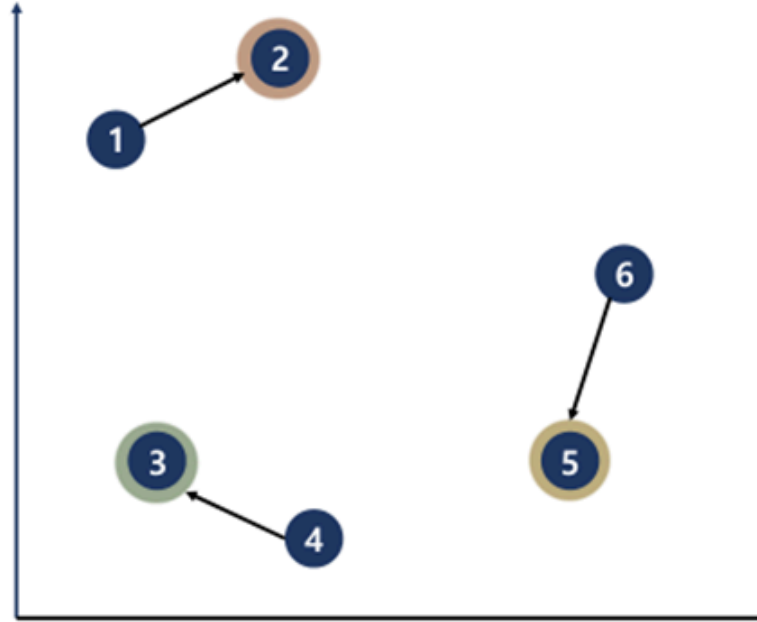
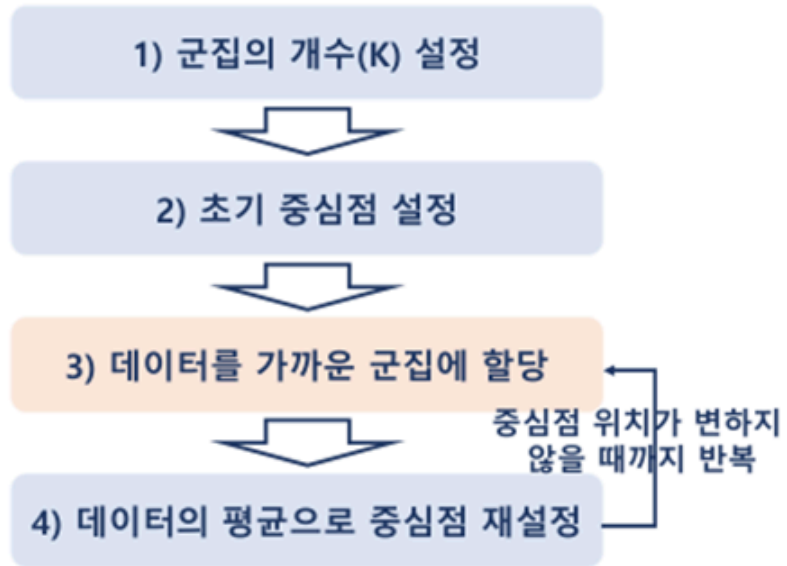
K-Means 군집화

각 데이터와 중심점의 거리를 측정 후 가장 가까운 그룹에 할당하여 K개의 군집으로 묶는 방법



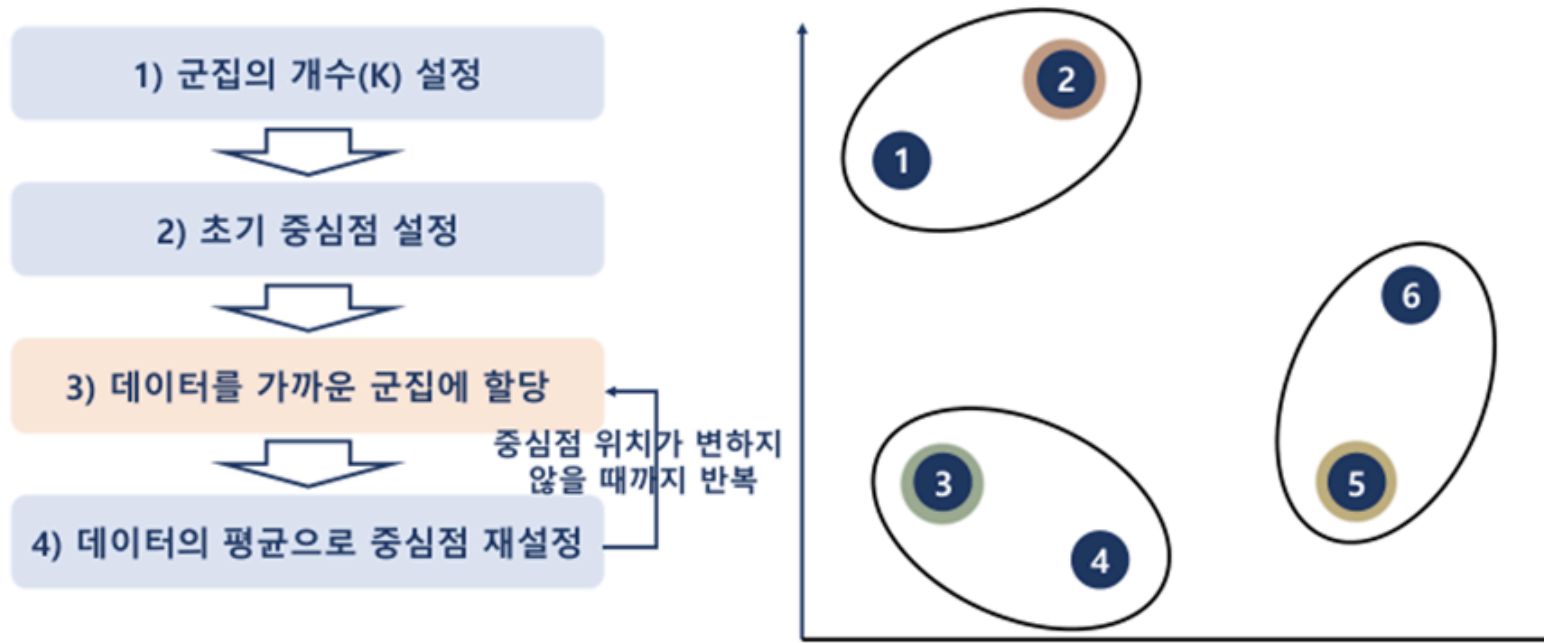
K-Means 군집화

각 데이터와 중심점의 거리를 측정 후 가장 가까운 그룹에 할당하여 K개의 군집으로 묶는 방법



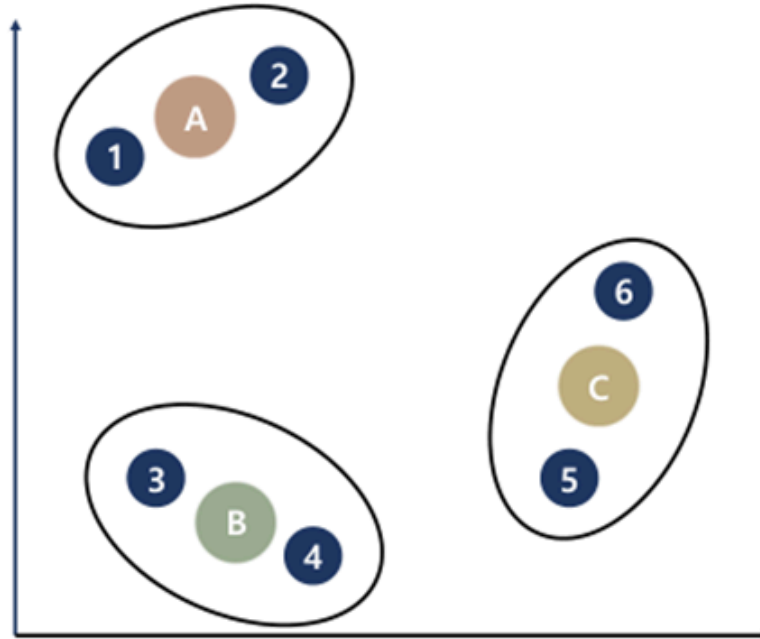
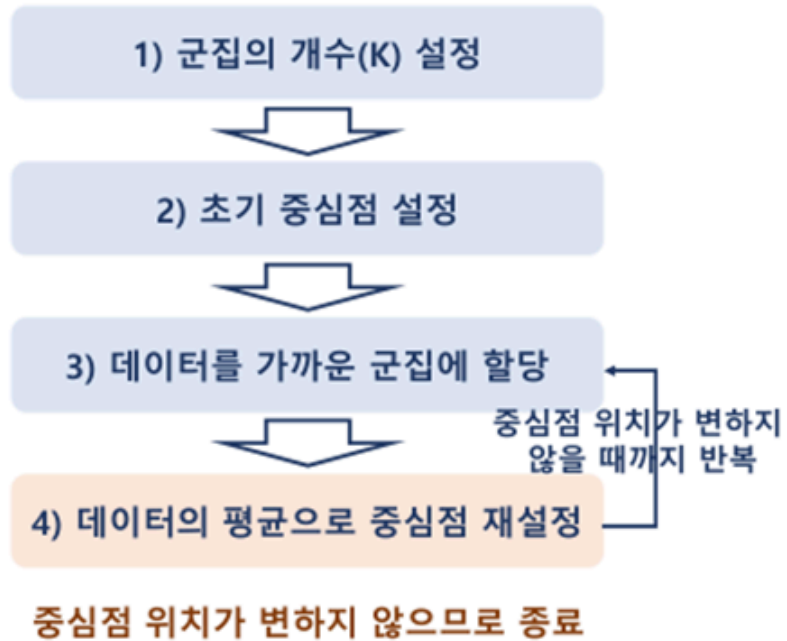
K-Means 군집화

각 데이터와 중심점의 거리를 측정 후 가장 가까운 그룹에 할당하여 K개의 군집으로 묶는 방법



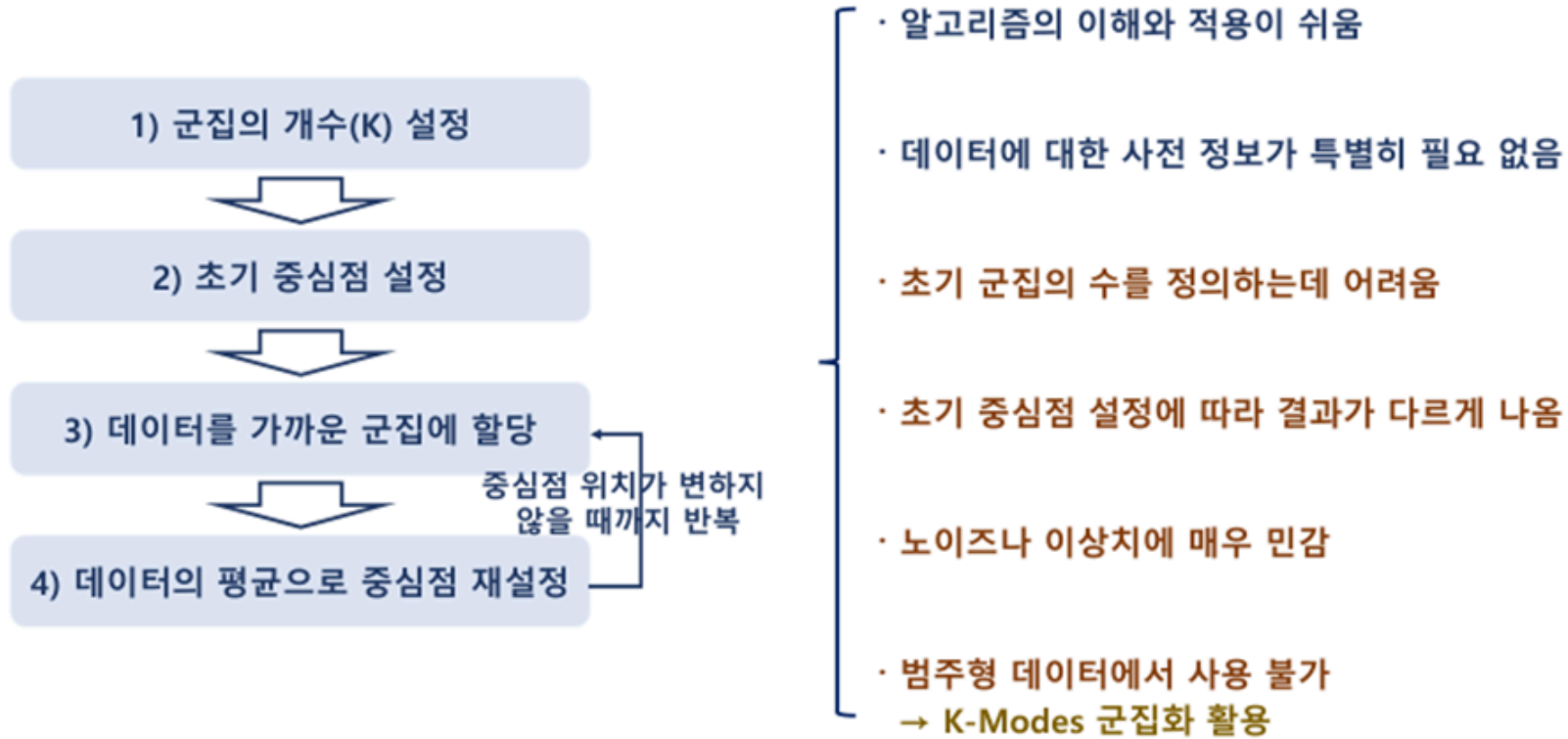
K-Means 군집화

각 데이터와 중심점의 거리를 측정 후 가장 가까운 그룹에 할당하여 K개의 군집으로 묶는 방법



K-Means 군집화

각 데이터와 중심점의 거리를 측정 후 가장 가까운 그룹에 할당하여 K개의 군집으로 묶는 방법



K-Means 군집화

최적의 군집 개수 K를 결정하는 방법?

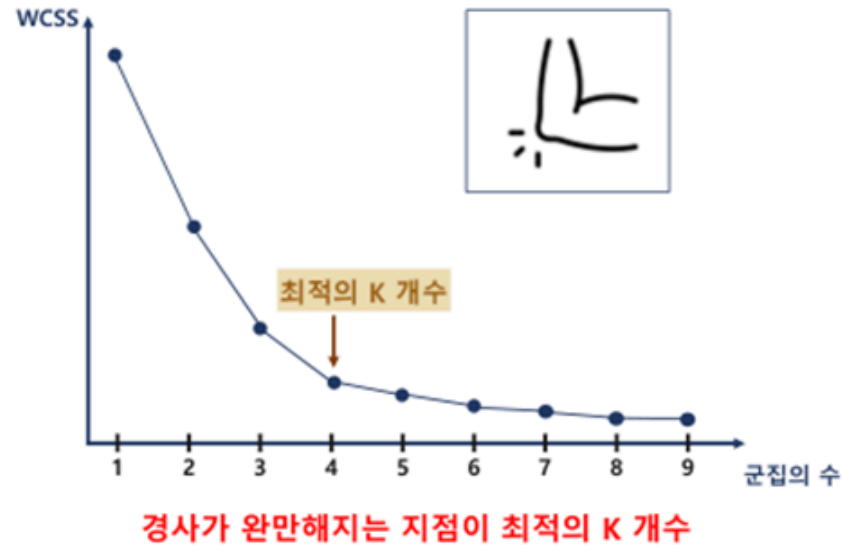
엘보우 기법(Elbow Method)

WCSS 값과 군집의 개수를 두고 비교 한 그래프를 통해 최적의 K 값을 선택하는 기법

$$WCSS = \sum_{C_k}^{C_m} \left(\sum_{d_i \in C_k}^{d_m} distance(d_i, C_k)^2 \right)$$

(Within Clusters Sum of Squares)

- C : 군집(Cluster)의 중심 값
- d : 클러스터 내에 있는 데이터



DBSCAN 군집화

데이터의 밀도를 기반으로 서로 가까운 데이터들을 군집으로 묶는 방법

- 데이터 밀도: 정보 작성의 기초가 되는 데이터에서 단위 길이, 단위 면적, 단위 부피마다 기억된 문자 수. (어휘 혼종어 정보·통신)

■ DBSCAN 군집화

데이터의 밀도를 기반으로 서로 가까운 데이터들을 군집으로 묶는 방법

1) 포인트 임의로 선택



2) Epsilon 거리 내 모든 데이터 탐색



3) Min Points 이상이면 군집에 할당



4) Min Points 이하이나 군집에 속한
포인트는 Border Point



5) 어느 군집에도 속하지 않는
포인트는 이상치



DBSCAN 군집화

데이터의 밀도를 기반으로 서로 가까운 데이터들을 군집으로 묶는 방법

1) 포인트 임의로 선택



2) Epsilon 거리 내 모든 데이터 탐색



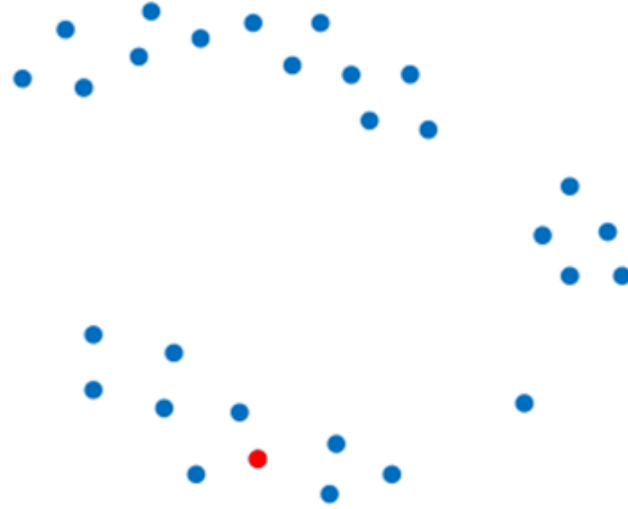
3) Min Points 이상이면 군집에 할당



4) Min Points 이하이나 군집에 속한
포인트는 Border Point



5) 어느 군집에도 속하지 않는
포인트는 이상치



DBSCAN 군집화

데이터의 밀도를 기반으로 서로 가까운 데이터들을 군집으로 묶는 방법

1) 포인트 임의로 선택



2) Epsilon 거리 내 모든 데이터 탐색



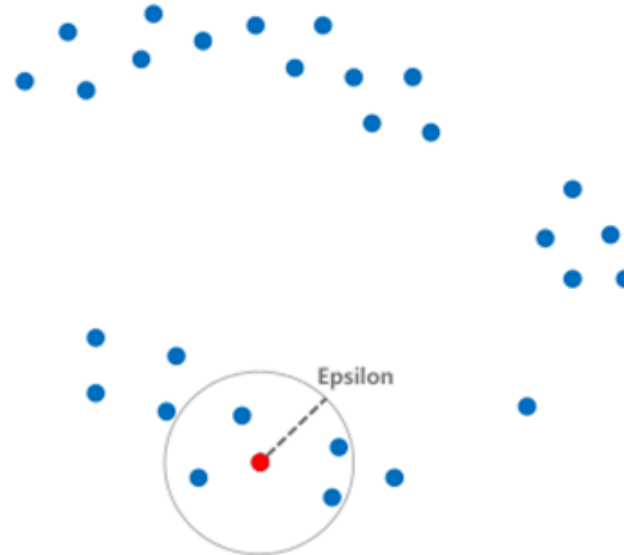
3) Min Points 이상이면 군집에 할당



4) Min Points 이하이나 군집에 속한
포인트는 Border Point

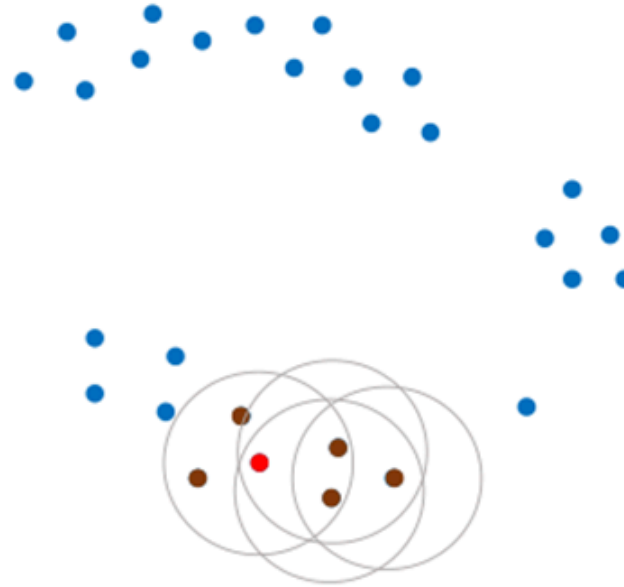
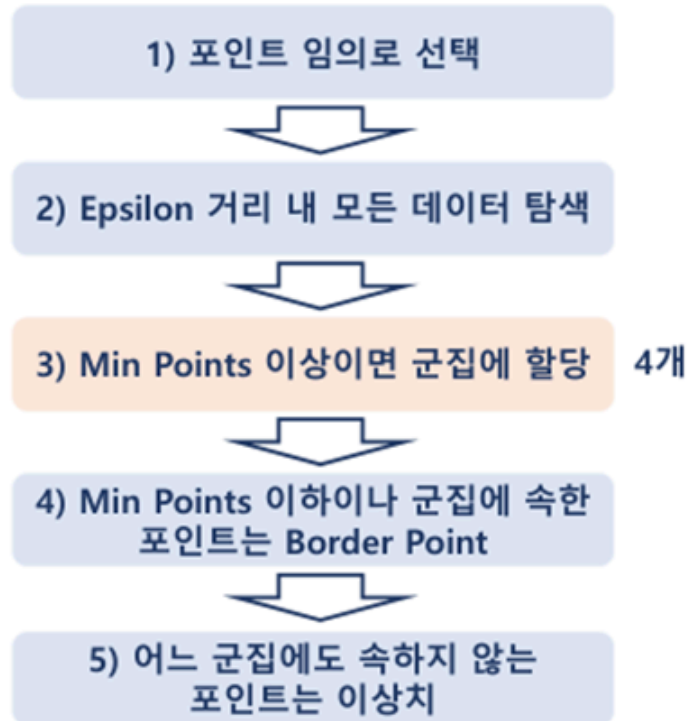


5) 어느 군집에도 속하지 않는
포인트는 이상치



DBSCAN 군집화

데이터의 밀도를 기반으로 서로 가까운 데이터들을 군집으로 묶는 방법



DBSCAN 군집화

데이터의 밀도를 기반으로 서로 가까운 데이터들을 군집으로 묶는 방법

1) 포인트 임의로 선택



2) Epsilon 거리 내 모든 데이터 탐색



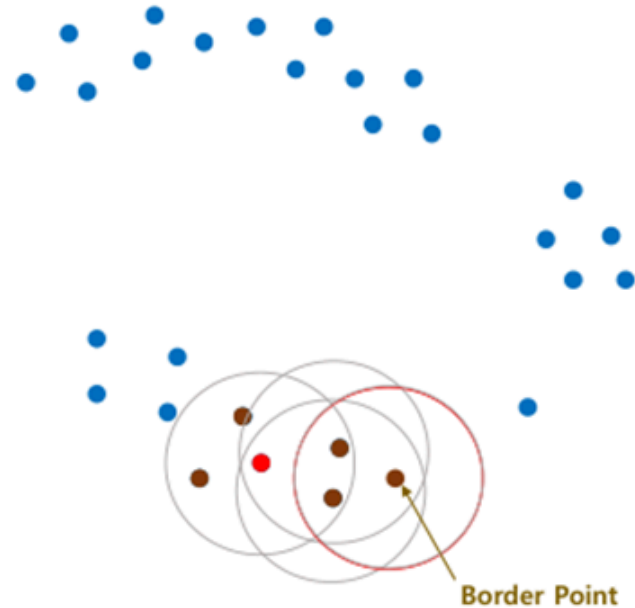
3) Min Points 이상이면 군집에 할당



4) Min Points 이하이나 군집에 속한
포인트는 Border Point



5) 어느 군집에도 속하지 않는
포인트는 이상치



■ DBSCAN 군집화

데이터의 밀도를 기반으로 서로 가까운 데이터들을 군집으로 묶는 방법

1) 포인트 임의로 선택



2) Epsilon 거리 내 모든 데이터 탐색



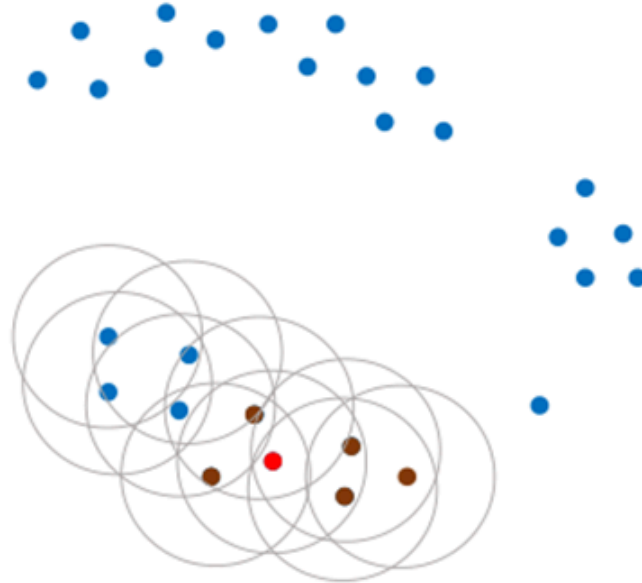
3) Min Points 이상이면 군집에 할당



4) Min Points 이하이나 군집에 속한
포인트는 Border Point



5) 어느 군집에도 속하지 않는
포인트는 이상치



DBSCAN 군집화

데이터의 밀도를 기반으로 서로 가까운 데이터들을 군집으로 묶는 방법

1) 포인트 임의로 선택



2) Epsilon 거리 내 모든 데이터 탐색



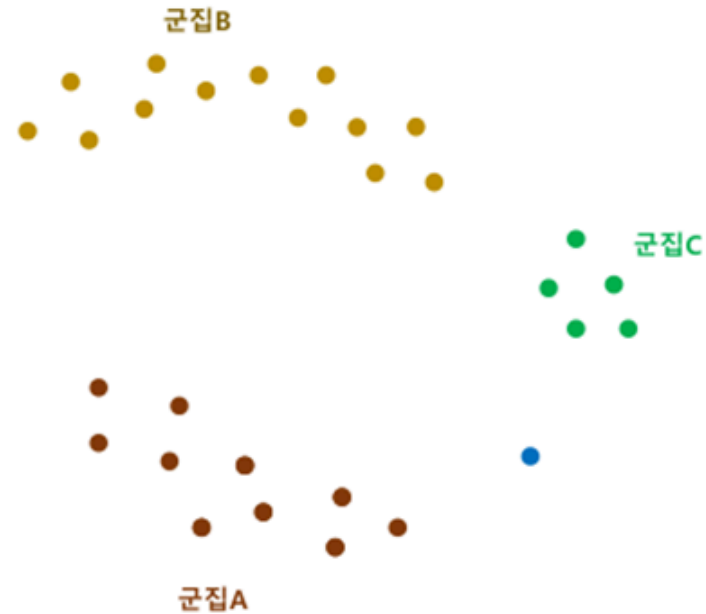
3) Min Points 이상이면 군집에 할당



4) Min Points 이하이나 군집에 속한
포인트는 Border Point



5) 어느 군집에도 속하지 않는
포인트는 이상치



■ DBSCAN 군집화

데이터의 밀도를 기반으로 서로 가까운 데이터들을 군집으로 묶는 방법

1) 포인트 임의로 선택



2) Epsilon 거리 내 모든 데이터 탐색



3) Min Points 이상이면 군집에 해당



■ K-Means와 DBSCAN 비교