

## 07장. 군집화

### 01. K-평균 알고리즘 이해

- 군집화 (clustering)에서 가장 일반적인게..
- 특정한 양의 점을 선택, 해당 중심에 가장 가까운 포인트를 선택. ) 이 과정을 반복.
- 선택된 포인트의 평균 지점을 중심점이 이동
- 중심점의 이동이 없을 때까지 위의 과정을 반복.

#### · K-평균의 장점

- 일반적 군집화에 가장 많이 활용.
- 쉽고 간단.

#### · K-평균의 단점

- 개수가 많으면 정확도 ↓
- 몇개의 군집을 선택할지 어려움.
- 반복횟수가 높으면 느려진다.

#### · 사이킷런 kmeans 파라미터.

- n\_clusters: 군집화할 개수 (평균 중심점의 개수)
- Init: 초기에 평균 중심점의 좌표를 선택할 방식.
- max\_iter: 최대 반복 횟수

#### · 군집화 알고리즘 테스트를 위한 데이터 생성

##### · make\_blobs()

→ 세벌 군집의 중심점과 표준 편차 개수  
가변이 큼.

##### · make\_classification()

→ 노이즈를 포함한 데이터 생성에 유용.

군집화용 데이터 생성기

### make\_blobs() 파라미터

- n\_samples: 생성할 총 데이터 개수
- n\_features: 데이터 개수
- centers: Init 값.
- cluster\_std: 생성된 군집 데이터의 표준 편차.

### 02. 군집평가

- 실용적 분석. (각 군집 간의 거리가 얼마나 효율적으로 분리돼 있는지..)

→ 동일 군집끼리만 잘 응집되고, 다른 군집끼리는 떨어져 있음.

실용적 계수를 통해 표현.

$$\text{실용적 계수 } S(i) = \frac{(b(i) - a(i))}{\max(a(i), b(i))}$$

· 1과 1 사이의 값.

· 1로 가까워질수록 군치 군집화 멀리, 0에 가까울수록 가까워진다.

· 전체 실루엣 계수의 평균값 :  $\text{silhouette\_score}()$  는 0~1의 값, 1에 가까울 수록 좋다.

### 03. 평균이동 (mean shift)

· K평균과 유사하게 군집의 중심으로 반복적으로 움직이면서 군집화를 수행.

· 데이터의 분포도를 이용해 군집 중심점을 찾음.

데이터 포인트가 모여있는 곳. (복잡한 밀도 함수를 이용) KDE (kernel density estimation) 이용.

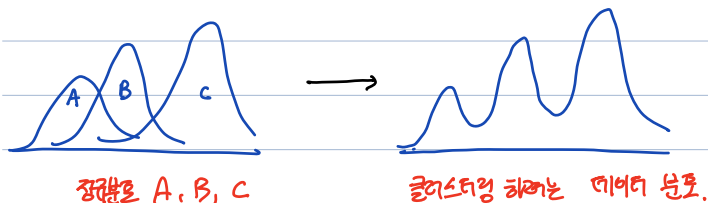
· 관측된 데이터 각각에 kernel을 적용하여 합한다.

· 데이터 간격을 나눠 kernel도 함수를 찾음. → 비모수적 가우시안 분포 함수 사용.

대역폭  $h$ : smoothing 에 이용.  $h$ 가 크면: 적은 수의 군집 중심  
 $h$ 가 작으면: 많은 수의 "

### 04. GMM (Gaussian Mixture Model)

· 군집화를 적용하면과 같은 데이터가 여러 개의 가우시안 분포를 가진 데이터 집합들이 섞여서 생성된 것이라는 가정 하..



여러 데이터 set → 이를 구성하는 여러개의 정규분포 곡선들을. → 개별 데이터가 이 중 어떤 정규 분포에 속하는지.

문득문득 ① 개별 정규 분포의 평균 / 분산

② 각 데이터가 어떤 정규분포에 해당하는지..

\* GMM은 kmeans 보다 유연하게 다양한 세트에 잘 적용될 수 있다는 장점. but 시간 ↑  
다원형은 많음

### 05. DBSCAN

· 밀도 기반 군집화의 대표적인 알고리즘.

#### 파라미터

· 임신을 위한 명목 (epsilon) : 개별 데이터를 중심으로 임신을

· 최소 데이터 개수 (min points) : 개별 데이터의 임신을 위한 명목에 포함되는 다 데이터의 개수.