

계층적 clustering (Hierarchy)

↳ 하나의 Data sample을 하나의 cluster로 보고 가장 유사도가 높은 cluster를 합치면서 cluster 갯수를 줄여나가는 방법.

$$k=N \Rightarrow k=N-1 \Rightarrow \dots \Rightarrow k=1$$

Hierarchy clustering에는 ① 비귀납적 방법 ② 귀납적 방법이 있다.

① 비귀납적 방법

i) Centroid : 두 cluster의 중심점 (centroid)를 정의한 다음, 두 중심점의 거리를 cluster 간의 거리로 정의한다.

$$d(u, v) = \|c_u - c_v\|_2 \quad (\text{여기서 } c_u \text{와 } c_v \text{는 각각 두 cluster } u \text{와 } v \text{의 중심점})$$

ii) Single : 클러스터 u의 모든 데이터 i와 클러스터 v의 모든 데이터 j의 모든 조합에 대한 거리를 측정해서 최소값을 구한다. 최근 거리 (Nearest point) 방법이라고도 한다.

$$d(u, v) = \min(\text{dist}(u[i], v[j]))$$

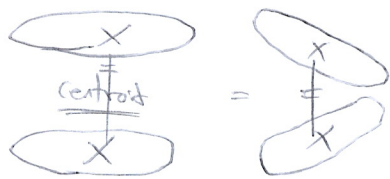
iii) Complete : 클러스터 u의 모든 데이터 i와 클러스터 v의 모든 데이터 j의 모든 조합에 대한 거리를 측정해서 가장 큰 값을 구한다. Farthest point Algorithm 또는 Vor Hees Algorithm이라고도 한다.

$$d(u, v) = \max(\text{dist}(u[i], v[j]))$$

iv) Average : 클러스터 u의 모든 데이터 i와 클러스터 v의 모든 데이터 j의 모든 조합에 대한 거리를 측정후, 평균을 구한다. |u|와 |v|는 각각 두 클러스터의 원소의 갯수를 뜻한다.

$$d(u, v) = \frac{\sum_{i,j} \text{dist}(u[i], v[j])}{|u||v|}$$

Diagram



⇒ Centroid로 클러스터 간의 거리를 정의한다면 원소 그림과 같이 클러스터 간의 거리는 같지만, 진짜 똑같다고 보기 힘들. 따라서, Single, Complete, Average가 사용.

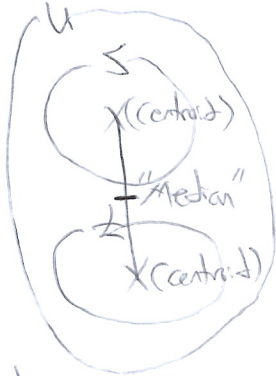


⇒ Single 방식: 가장 짧은 거리 방법
Complete 방식: 가장 긴 거리 방법.
Average 방식: 평균 거리 방법.

✗ 비귀납적 방법의 문제점 : 계산량이 너무 많기 때문에, 많은 Data를 쓸 때 사용하기는 현실적으로 힘들. 따라서, 보통 귀납적 방법을 사용.

② 귀납적 방법

i) Median : 이 방법은 Agglomerative clustering에서 사용할 수 있는 귀납적 방법으로 Centroid 방법의 변종이다. 만약 cluster가 cluster S와 cluster T가 결합하여 생겼다면, cluster U의 중심점은 새로 계산하지 않고 원래 클러스터의 중심점의 중간점을 사용한다.



ii) Weighted : 이 방법도 Agglomerative clustering에서 사용할 수 있는 귀납적 방법이다. 만약 cluster U가 cluster S와 cluster T가 결합하여 생겼다면, 다음과 같이 원래 cluster까지의 즉 거리의 가중치를 사용한다.

iii) Ward : 이 방법도 Agglomerative clustering에서 사용할 수 있는 귀납적 방법이다. 만약 cluster U가 cluster S와 cluster T가 결합한 것이라면, 두 cluster간의 거리의 가중치에서 원래의 즉 클러스터 사이의 거리를 보정한 값을 사용한다.

$$d(u, v) = \sqrt{\frac{|u|+|s|}{|u|+|s|+|t|} d(v, s)^2 + \frac{|u|+|t|}{|u|+|s|+|t|} d(v, t)^2 - \frac{|u| \cdot d(s, t)}{|u|+|s|+|t|}}$$

여기서 |·|는 클러스터의 원소 개수를 의미함.

- Scipy의 Hierarchy clustering을 쓸 수 있음.

$Z = \text{linkage}(X, 'ward')$ # X가 100

array([[52.,
Index #
[14.,

53.
Index #
79.

0.04151,
두 개의 cluster간의 거리
(Paranet의 거리)
0.05919,

2.
원래의 개수
150번 cluster
150번 cluster

[62.,
152.

0.1726, 3.]

(([33,
[60,
[62])

"가장 처음의 개수"

↑

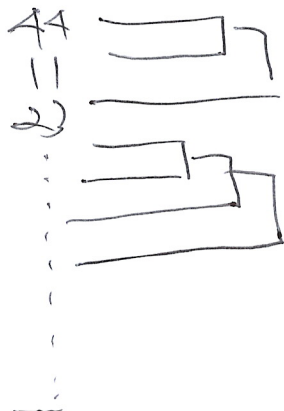
150번 cluster

←

150번 cluster

←

- Dendro gram



→ 사실 순서: 비슷한 Data끼리 외도록 재배열된 것임 (원래는 1, 2, 3, ..., N)

→ 일종의 "Quasi Sorting"이 됨.

가우시안 혼합 모형과 EM 방법

가우시안 혼합 모형 (Gaussian Mixture)

→ 실수 값을 출력하는 확률 변수 X 가 사실 z 에 속하지 않는 (관측되지 않는) K -class 카테고리 확률 변수 z 의 값에 따라, 다른 가설변과 분산을 가지는 복수의 가우시안 정규 분포를 띄워 띄워진 모형을 가우시안 혼합 모형이라고 한다.

$$p(x) = \sum_{k=1}^K p(z) \cdot p(x|z) = \sum_{k=1}^K \theta_k N(x | \mu_k, \Sigma_k)$$

- $p(x)$: 전체 분포, $p(x|z)$: 혼합 모형의 각 성분이 되는 개별적인 연속 확률 분포
- $p(z)$: 카테고리 확률 분포, θ_k : 카테고리 확률 분포의 요소. 모든 성분중에 특정한 z 가 선택될 확률.
- $\sum \theta_k = 1$: 카테고리 확률 분포의 요소 제한 조건.

베르누이-가우시안 혼합 모형 (Bern - Gaussian Mixture)

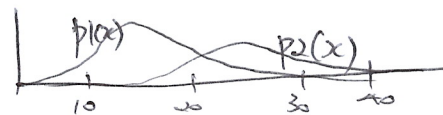
- 카테고리가 2개인 가우시안 혼합 모형은 베르누이-가우시안 혼합 모형이라고 한다.

active component

Component densities

observations

comp 1 comp 2
 $\theta_1 = 0.15$ $\theta_2 = 0.25$



14.2

20.9

Bern - Gaussian Mixture
 이 문제점은 우리가 아는 것
 observation #이와, 우리가
 원하는 것은 이 확률 모형에서
 나타나는 가설의 확률.
 ↓
 그것을 보거나 각각의
 요소 확률로 알고 싶은 것임
 ↓
 우리는 전체 $p(x)$ 는 아니지만
 개별 성분 확률 모형을
 우리가 어떻게 해서 모아야 하는 것임

EM (Expectation - Maximization)

- 혼합 모형의 모수추정에서 중요한 역할을 하는 것중에 하나가 바로 각 data가 어떤 category에 속하는지를 알려주는 소위 responsibility이다. 이 값을 responsibility라고 한다.

$$r_{ik} = p(z_i = k | x_i) = \frac{p(z_i = k) \cdot p(x_i | z_i = k)}{p(x_i)}$$

$$= \frac{p(z_i = k) \cdot p(x_i | z_i = k)}{\sum_{k=1}^K p(x_i, z_i = k)}$$

$$= \frac{p(z_i = k) \cdot p(x_i | z_i = k)}{\sum_{k=1}^K p(z_i = k) \cdot p(x_i | z_i = k)}$$

- 가설인 혼합모형의 경우 다음과 같이 정리된다.

$$r_{ik} = \frac{\theta_k N(x_i | \mu_k, \Sigma_k)}{\sum_{k=1}^K \theta_k N(x_i | \mu_k, \Sigma_k)}$$

\Rightarrow 이 때 모수추정 responsibility를 측정한다.
 $(\theta_k, \mu_k, \Sigma_k) \Rightarrow r_{ik}$

* r_{ik} 은 i번째 data x_i 가 cluster k에 속할 확률을 나타낸다.

- 일반적으로 Log Likelihood 함수를 미분하여 0이 되도록 하는 모수를 찾는 과정을 반복한다.

$$0 = - \sum_{i=1}^N \frac{p(z_i = k) p(x_i | z_i = k)}{\sum_{k=1}^K p(z_i = k) \cdot p(x_i | z_i = k)} \sum_k (x_i - \mu_k)$$

- 이를 정리하면 $\sum_{i=1}^N r_{ik} (x_i - \mu_k) = 0$, $\mu_k = \frac{1}{N_k} \sum_{i=1}^N r_{ik} x_i$

$$\Rightarrow N_k = \sum_{i=1}^N r_{ik}$$

- k cluster에 속하는 데이터의 수와 비슷한 의미를 가진다. 즉 μ_k 는 k cluster에 속하는 데이터의 샘플평균과 같은 의미이다.

- 마찬가지로 log likelihood를 \sum_k 로 미분하여 cluster k에 속할 확률을 구하는 과정을 반복한다.

$$\sum_k = \frac{1}{N_k} \sum_{i=1}^N r_{ik} (x_i - \mu_k)(x_i - \mu_k)^T$$

- 마지막으로 log likelihood를 θ_k 로 미분하여 cluster k에 속할 확률을 구해야 하는데, 이 때 cluster k의 모수가 가지는 제한조건으로 인해 Lagrange Multiplier를 추가해야 한다.

$$L + \lambda \left(\sum_{k=1}^K \theta_k - 1 \right)$$

이를 θ_k 로 대체하여 0이 되는 것은 없으면 $\theta_k = \frac{N_k}{N}$

- 이 세가지 수를 모두 responsibility로 부터 r_{ik} 를 구하고 있다.

$$r_{ik} \Rightarrow (\theta_k, \mu_k, \Sigma_k)$$

- 또한 방법은 r_{ik} 와 responsibility를 비교하여 클러스터링 정확도를 높여주는 방법.

- "E" step에서는 우리가 현재까지 알고있는 r_{ik} 가 정확하다고 가정하고 이를 사용하여 각 Data가 어느 카테고리에 속하는지, 즉 responsibility를 추정.

$$(\theta_k, \mu_k, \Sigma_k) \Rightarrow r_{ik}$$

- "M" step에서는 우리가 알고있는 responsibility가 정확하다고 가정하고 r_{ik} 를 추정.

$$r_{ik} \Rightarrow (\theta_k, \mu_k, \Sigma_k)$$

\Rightarrow 이를 반복하면, r_{ik} 와 responsibility를 동시에 점진적으로 개선 가능

Clustering

각각의 Data에 대해 responsibility를 알려주면 responsibility가 가장 큰 카테고리를 찾아내서 그 카테고리가 어떤 카테고리에 속하는지 알 수 있다. 즉, clustering 가능.

$$k_i = \arg \max_k r_{ik}$$