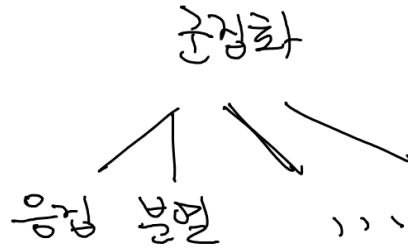


공간과 시간 (이점) 정리. 군집화의 일환으로 여러 특성을 가진 샘플들을 하나의 군집으로 묶는 비지도학습의 방법
 \Rightarrow Clustering

군집화를 할 때 샘플들을 군집화하는 기준.

거리나 유사도

1. 미노프스 거리 $\sum (x_1 - x_2)^{\frac{1}{p}}$
 $(p=2)$ 유클리드 $\sqrt{(x_1 - x_2)^2}$
 $(p=1)$ 맨하튼 $|x_1 - x_2|$



k-means

Clustering의 방법에서 응집, 분열 등 중의 분열 알고리즘이 속하는 방식

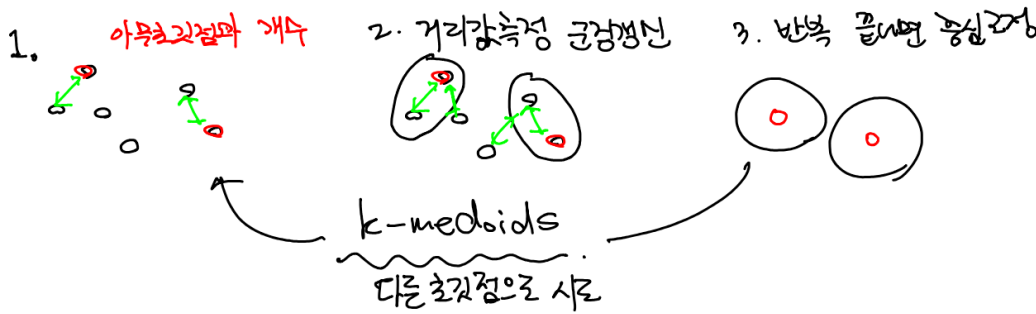
군집의 개수가 주어지며 Dmax, Dmin, Davg를 이용하여 군집의 중심을 부정하며 마지막이 이전과 군집의 중심이 같게 되면 끝나는 알고리즘

초깃값과 outlier이 민감하다 (거리측정에 대한 오차가 발생).

해결책 : 표준편차로 정렬값을 판단함에 따라 초깃값을 다변하여 반복 (k-medoids)

장점 : 타 알고리즘에 비해 빠르다.

알고리즘 순서



단점 : 초깃값과 outlier에 약하다.

SVM : $\sum \frac{1}{2} \frac{||x||^2}{n}$ 이렇게 분류시킬 것임에 대한 방법론

분류하여 양과 음의 클래스

Category Data : 특성이 있는 범주형 데이터

Data Embedding : 2차 데이터를 벡터형으로 변환

Decision Function 인 Hyperplane을 통해 나눈다.

$Wx + b > 1$ Positive
 $Wx + b < -1$ Negative

장점 : 분류에 예측이 틀다수 불가능하다.

좋은 Hyperplane은? 있어야 하는 Data가 나뉘는 Hyperplane

단점 : kernel 및 파라미터 값은 어디로 해야 할지 모르겠다.

Margin을 최대화 할 때 좋은 Hyperplane.
 (여백)

SVM은 두 Margin의 끝이 걸려있는 샘플의 벡터이다.

kernel Gamma : 분류가 어려운 데이터를 Gamma의 파라미터 크기를 통해 일정한 샘플의 위치를 조정하는 기술

Rbf kernel

Gamma의 Parameter를 찾는 기술 : Grid Search

$WTW + C = \text{Cost} \uparrow \text{margin} \downarrow / \text{Cost} \downarrow \text{margin} \uparrow$

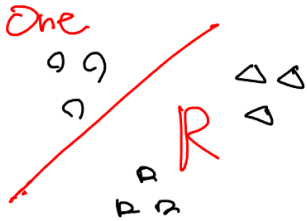
위키과정을 가진 SVM을 Soft SVM이라 부른다.

장점 : 실생활에 적용이 가능. 예측, 분류 모델에 잘라 사용이 가능하리.

단점 : Gamma. C 파라미터를 여러번 Test를 통해 찾아야한다.

Multi의 SVM

One vs Rest



클래스 각각의 이진 SVM을
호출한다.

One vs One



$\frac{k(k-1)}{2}$ 개 쌍의 SVM을

kernel 함수를 사용.

Out 보라 바래 있음.