

# (K-means clustering)

2-18/07/08

- 독립변수의 특성이 유사한 데이터 그룹  $\Rightarrow$  클러스터 (cluster)

- 주어진 데이터를 여러개의 클러스터로 구분하는 것  $\Rightarrow$  클러스터링 (clustering)

## K-means clustering

- 가장 간단하고 빠른 cluster Algorithm 중 하나.

- 독립변수 값이 유사할수록 데이터가 클러스터의 중심 (centroid)  $\mu_k$ 과  
각 데이터가 소속될 cluster를 반복해서 찾는 것. 이 최소화된 독립변수 값을 inertia라 한다.

$$J = \sum_{k=1}^K \sum_{i \in C_k} d(x_i, \mu_k) \quad d(x_i, \mu_k) = \|x_i - \mu_k\|^2$$

- 새 Algorithm은 다음과 같다.

- ① 임의의 중심점  $\mu_k$ 를 고른다.
- ② 중심에서 각 sample data 가와의 거리를 계산.
- ③ 각 데이터 샘플에서 가장 가까운 중심을 선택하여 클러스터 생성.
- ④ 다시 만들어진 클러스터에 대해 중심을 다시 계산하여 1~4회를 반복한다.

## K-means++

- K-means algorithm은 최적의 중심점을 설정하기 위한 알고리즘.

- ① 중심점을 설정할 점의 집합  $S$  준비.
- ② 임의의 하나의 중심  $\mu_0$ 를 선택하여  $S$ 에 넣는다.
- ③  $S$ 에 속하지 않는 모든 샘플  $x_i$ 에 대해 거리  $d(S, x_i)$ 를 계산.  
 $d(S, x_i)$ 는  $S$ 안의 모든 점을  $\mu_k$ 에 대해  $d(\mu_k, x_i)$ 를 계산하여 가장 작은 값 선택
- ④  $d(S, x_i)$ 에 비례한 확률로 다음 중심  $\mu$ 를 선택.
- ⑤ K개의 중심을 선택할 때까지 반복.
- ⑥ K-means Algorithm 사용.