

# clustering

Sunday, September 19, 2021 5:51 PM

you have

Clustering refers to partitioning a set into "dissimilar" subsets of "similar" elements.

---

Usually a well-defined objective.

E.g. k-means, k-median, k-center, diameter.  
— define

or with the goal of recovering some ground truth.

---

— all are NP-hard

---

No universal clustering criterion.



Depends on the context/application.

---

K-center: - start with any point in given set  
as first center.  $C = \{C_1\}$

Repeat [ - add farthest point to  $C$   
K-1 times

---

Thm. Greedy algorithm is a factor 2 approximation.

pf. Suppose OPT is  $R$ .

Claim: For the centers  $C_1, \dots, C_K$  found by GREEDY,  $\max$  distance to nearest center  $\leq 2R$ .

If not,

$\Rightarrow \exists K+1$  pts  $C_1, \dots, C_{K+1}$

s.t.  $d(C_i, C_j) > 2R$ .

$\Rightarrow$  No two of  $\{C_1, \dots, C_{K+1}\}$  can belong  
to same  $d \leq R$ .

$\Rightarrow$  No ms of  $\{c_1, \dots, c_{k+1}\}$   
to same cluster of radius  $r$ .

## The Spectral Approach

- Project to span of top  $k$  singular vectors of  $A \in \mathbb{R}^{n \times d}$
- Cluster in  $\mathbb{R}^k$ .

Idea: this should shrink distance between  $x$  and nearest center.

$$d\left( n \begin{pmatrix} A \\ \vdots \end{pmatrix} - \begin{pmatrix} c_1 \\ c_2 \\ \vdots \\ c \end{pmatrix} \right)$$

centers.

$$\tilde{\sigma}(C) = \frac{\|A - C\|_2}{n} : \text{average variance of clusters.}$$

## Clusters :

Thm. If  $\|C_i - C_j\| > \frac{15}{\epsilon} \cdot \sigma(C)$  & if  $\epsilon$  and each cluster has  $\leq n$ , then Spectral Clustering finds  $C'$  that differs from  $C$  in at most  $\epsilon^2 \cdot n$  points.

Algo. 1. Project to top  $k$  right Singular vectors of  $A$ .

Repeat  $k$  times. 2. Take a random row, include all points within distance  $\frac{6K\sigma(C)}{\epsilon}$ .

$$A_k = \underset{D: \text{rank}(D) \leq k}{\operatorname{argmin}} \|D - A\|_2$$

Lem. for any  $C$  of rank  $K$ ,  $\|A_k - C\|_F^2 \leq 8k \|A - C\|_2^2$  any  $A$ .

Pf. (\*)  $\|A_k - C\|_F^2 \leq 2k \|A_k - C\|_2^2$

Since  $A_k - C$  has rank  $\leq 2k$ .

Since  $A_K - C$  has  $\text{rank} \leq 2K$ .

$$\|A_K - C\|_2 \leq \|A_K - A\|_2 + \|A - C\|_2 \\ (\star\star) \quad \leq 2\|A - C\|_2$$

$$(\star) + (\star\star) \Rightarrow \|A_K - C\|_F^2 \leq 8K \|A - C\|_2^2.$$

---

Pf (of Thm). Let  $v_i$  be  $i^{\text{th}}$  row of  $A_K$ .

Claim. Most  $v_i$  are within distance  $\frac{3K\sigma(C)}{\varepsilon}$  of their center.

$$\text{Let } B = \left\{ i : \|v_i - C\| > \frac{3K}{\varepsilon} \sigma(C) \right\}.$$

$$\text{Then } \|A_K - C\|_F^2 \geq |B| \cdot \frac{9K^2}{\varepsilon^2} \cdot \sigma(C)^2 \\ \leq 8K \sigma(C)^2 \cdot n$$

$$\Rightarrow |B| < \frac{\varepsilon^2}{K} \cdot n.$$

For  $i, j \in$  same cluster and  $\notin B$ ,

...  $v_i - v_j$  ...  $\rightsquigarrow$



$$\|v_i - v_j\| \leq \frac{6K}{\epsilon} \sigma(C).$$

✓

$i, j$  different centers,  $\notin B$

$$\|v_i - v_j\| > \frac{15K}{\epsilon} \sigma(C) - \frac{6K}{\epsilon} \sigma(C) = \frac{9K}{\epsilon} \sigma(C).$$

Hence if we pick point not in  $B$  as the seed, all  $\leftarrow$  time, all points not in  $B$  will be correctly classified.

$$\begin{aligned} P_1(\text{we pick point } \notin B) &\geq \left(1 - \frac{\epsilon^2}{K}\right) \cdot \left(1 - \frac{(\epsilon - \epsilon^2)^2}{K}\right)^{*1} \\ &\geq 1 - \frac{\epsilon}{K} \cdot K = 1 - \epsilon. \end{aligned}$$


---

Example Mixture of  $K$  Gaussians, each with max covariance  $\sigma^2$ .

$$\text{Then } r(C) \leq \sigma \sqrt{K} \quad (\text{distance to zero center})$$

$$\text{Separation needed: } 15K \frac{\sigma(C)}{\epsilon} \leq 15K \frac{\sigma^2}{\epsilon} \cdot \sigma.$$

Example 2 Stochastic Block Models.

Example 2

Stochastic Block Models .  
or Planted Partitions .