

Clustering

Sunday, September 19, 2021 5:51 PM

challenge

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

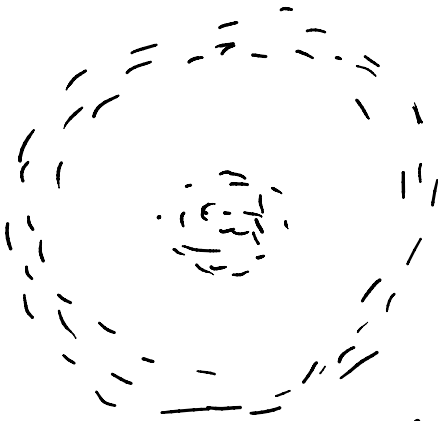
Usually a well-defined objective.

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

or with the ^{define} 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 $\left[\begin{array}{l} \text{- add farthest point to } C \\ \text{K-1 times} \end{array} \right.$

Thm. Greedy algorithm is a factor 2 approximation.

Pr. Suppose OPT is γ .

Claim: For the centers c_1, \dots, c_k found by
GREEDY, max distance to nearest center $\leq 2\gamma$.

$\Rightarrow \exists k+1$ pts c_1, \dots, c_{k+1}

st. $d(c_i, c_j) > 2\gamma$.

\Rightarrow No two of $\{c_1, \dots, c_{k+1}\}$ can belong
to same cluster of radius γ .

The Spectral Approach

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

- Cluster in \mathbb{R}^k .

... between

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

$$n \begin{pmatrix} A \end{pmatrix} - \begin{pmatrix} \text{--- } c_1 \text{ ---} \\ \text{--- } c_2 \text{ ---} \\ C \end{pmatrix}$$

$$\sigma^2(C) = \frac{\|A - C\|_2}{n} : \text{average variance of clusters.}$$

Thm. If $\|C_i - C_j\| > 15 \frac{k}{\epsilon} \cdot \sigma(C) \forall i \neq j$ and each cluster has $\leq \epsilon n$, 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. $\left[\begin{array}{l} 2. \text{ Take a random row, include all} \\ \text{points within distance } \frac{6k\sigma(C)}{\epsilon}. \end{array} \right.$

$$\|D - A\|.$$

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

Lemma. 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 $\leq 2k$.

$$\|A_k - C\|_2 \leq \|A_k - A\|_2 + \|A - C\|_2$$

(**) $\leq 2\|A - C\|_2$

(*) + (**) $\Rightarrow \|A_k - C\|_F^2 \leq 8k \|A - C\|_2^2$.

PF (of Thm). Let v_i be i^{th} row of A_k .

claim. Most v_i are within distance $\frac{3k\sigma(C)}{\epsilon}$
of their center.

Let $B = \{i : \|v_i - C\| > \frac{3k\sigma(C)}{\epsilon}\}$.

Then $\|A_k - C\|_F^2 \geq |B| \cdot \frac{9k^2}{\epsilon^2} \cdot \sigma(C)^2$

$\leq 8k\sigma(C)^2 \cdot n$

$$\Rightarrow |B| < \frac{\epsilon^2}{k} \cdot n.$$

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

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

i, j different clusters, $\notin B$

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

Hence if we pick point not in B as the seed, all k times, all points not in B will be correctly classified.

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

Example 1 Mixture of k Gaussians, each with max covariance σ^2 .

Then $r(C) \leq C_1 \sigma$ (max distance to our center in one direction)

Separation needed: $\frac{15k}{\epsilon} \sigma(C) = O\left(\frac{k}{\epsilon} \cdot \sigma\right)$ between centers

(a bit worse than what we got)

Example 2 Stochastic Block Models.
or Planted Partitions.

or Planted Partitions.

p	q	q
	p	
		p

$i, j \in$ same block \leftarrow
 $R((i, j) \in E) = \begin{cases} p \\ q \end{cases}$ different blocks.

A : adjacency matrix of G .

$$E(A) = C = \begin{pmatrix} p \dots p & | & q \dots q \\ \vdots & & \vdots \end{pmatrix}$$

Problem: Reconstruct "planted" partition.

Apply spectral algorithm.

$$\| \mathbf{1}_i - \mathbf{1}_j \|^2 = (p-q)^2 \cdot \frac{n}{k}$$

different clusters.

What about $\|A - C\|_2$?

$A - C$ is a random matrix with $E(A - C) = 0$ and independent entries.

Then, R random with independent entries $\in [-1, 1]$
 $E(R_{ij}) = 0$, $E(R_{ij}^2) \leq \sigma^2$, $\|R\| \in (2 + o(1)) \sigma \sqrt{n}$.

$\mathbb{E}(R_{ij}) = 0$ $\mathbb{E}(R_{ij}^2) \leq \sigma^2$.
 Then with prob. $1 - o(1)$, $\|R\|_2 \leq (2 + o(1)) \sigma \sqrt{n}$.

So, $\sigma(C)^2 = \frac{\|A - C\|_2^2}{n} = O(p)$

So by the spectral clustering theorem,
 it suffices to have

$$\|r_i - r_j\|^2 \geq (p - q)^2 \cdot \frac{n}{k} > c \frac{k^2}{\varepsilon^2} p > \left(\frac{15k}{\varepsilon}\right)^2 \cdot \sigma(C)^2$$

$$\text{or } |p - q| > c' \frac{k^{3/2}}{\varepsilon} \sqrt{\frac{p}{n}}$$

if we set $p = \frac{a}{n}$, $q = \frac{b}{n}$

then this says. $|a - b| = \Omega(k^{3/2}) \cdot \sqrt{a}$

This is information theoretically tight up to
 constant factors. suffices.

In fact for $k \geq 2$, $(a - b)^2 > 2a$
 is necessary and sufficient!

or uniform matrix bound?

Q. How to find the random matrix bond?

Planted Clique.