

clustering

Sunday, September 19, 2021 5:51 PM

you have

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

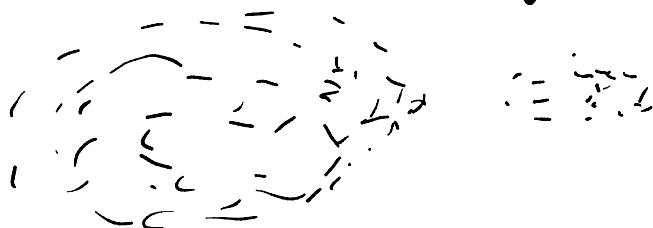
Usually a well-defined objective.

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

or with the goal of recovering some ground truth.

- all are NP-hard

No universal clustering criterion.



Depends on the context/application.

not

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 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 .

between

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

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

centers.

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

Thm. If $\|c_i - c_j\| > 15 \frac{k}{\epsilon} \cdot \sigma(C)$ $\forall i \neq j$ and each cluster has $\leq n$, 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 $6 \frac{k \sigma(C)}{\epsilon}$.

$$\|D - \Delta\|.$$

$$A_K = \underset{\mathcal{D}: \text{rk}(\mathcal{D}) \leq K}{\operatorname{arg\min}} \|\mathcal{D} - A\|_2$$

Lemma. For any C of rank K , $\|A_K - C\|_F^2 \leq 8K \|A - C\|_2^2$ for any A .

Pf. (*) $\|A_K - C\|_F^2 \leq 2K \|A_K - C\|_2^2$

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

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

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

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

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

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

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

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

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

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

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

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



i, j different clusters, $\notin B$

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

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

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

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

Then $r(C) \leq C_1 \sigma$ ($\max_{\text{center}} \text{distance to center in one direction}$)

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

(a bit worse than what we got)

Example 2 Stochastic Block Models.
or Planted Partitions.

on Planted Partitions.

p	q	q
	p	
		p

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

A: adjacency matrix of G.

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

Problem: Recover "planted" partition.

Apply spectral algorithm.

$$\|v_i - v_j\|^2 = .(p-q)^2 \cdot n$$

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

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

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

So by the spectral clustering theorem,
 it suffices to have

$$\|r_i - r_j\|^2 \geq (p-a)^2 \cdot n > c \frac{k^2}{\varepsilon^2} p > \left(\frac{5K}{\varepsilon}\right)^2 \cdot \sigma(C)^2$$

$$\text{or } |p-a| > c \frac{k}{\varepsilon} \sqrt{\frac{p}{n}}$$

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

$$\text{then this says. } |a-b| = \Omega(K) \cdot \sqrt{a}$$

suffices.

This is information theoretically tight up to constant factors.

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

... or random matrix bound?

Q. How to prove the random matrix bound?