

Coresets

Representative, smaller sample of a dataset

Coreset approximates whole dataset

Example: K-median problem

Given n points on plane $P \in \mathbb{R}^2$

Find K centers to serve the groups
 $c_1, \dots, c_k \in \mathbb{R}^2$

$$d(P, C) = \sum_{i=1}^n \min_{j \in [K]} \|P_i - c_j\|$$

"total distance to centers"

Why "median"? If $\|\cdot\| = L_1$,
 center of each group = median in
 each coordinate

$$\left(\text{k-means: } \sum_{i=1}^n \min_{j \in [K]} \|P_i - c_j\|_2^2 \right)$$

$$\text{OPT}(P, K) = \min_{|C|=K} d(P, C)$$

K-median problem: find C minimizing cost

NP-hard, so approximate:
Find C , $|C| \leq k$,

$$d(P, C) \leq (1+\epsilon) \cdot \text{OPT}(P, k)$$

Kollipour-Rao '99: $O\left(\left(\frac{1}{\epsilon}\right)^{O\left(\frac{1}{\epsilon}\right)} \cdot n \log n \log k\right)$

What about a streaming algorithm?

Today: generic technique for insertion streams
using coresets

A coreset for P is a weighted set of points S :

$$s_1, \dots, s_m \in [\Delta]^2$$

$$w_1, \dots, w_m \in \mathbb{Z}$$

$$\pi: [n] \rightarrow [m]$$

$$w_i = |\pi^{-1}(i)| = \#j \text{ s.t. } \pi(j) = i$$

$$\text{dist}(P, (S, w)) := \sum_{j \in [n]} \|P_j - S_{\pi(j)}\|$$

(k, ϵ) coreset if $\text{dist}(P, (S, w)) \leq \epsilon \cdot \text{OPT}(P, k)$

Let's us find appx. k -median solution:

$$d((S, w), C) := \sum_{j \in [m]} w_j \min_{i \in [k]} \|C_i - S_j\|$$

$$\tilde{C} := \underset{C}{\text{argmin}} d((S, w), C)$$

For any C ,

$$\begin{aligned}
 d(P, C) &= \sum_{i \in [n]} \min_j \|P_i - C_j\| \\
 &= \sum_{i \in [n]} \min_j (\|S_{\pi(i)} - C_j\| \pm \|P_i - S_{\pi(i)}\|) \\
 &= d((S, w), C) \pm \text{dist}(P, (S, w)) \\
 &= d((S, w), C) \pm \varepsilon \cdot \text{OPT}(P, k)
 \end{aligned}$$

Let $C^* = \arg \min_C d(P, C)$. Then

$$\begin{aligned}
 d(P, \tilde{C}) &= d((S, w), \tilde{C}) \pm \varepsilon \cdot \text{OPT}(P, k) \\
 d(P, C^*) &= d((S, w), C^*) \pm \varepsilon \cdot \text{OPT}(P, k)
 \end{aligned}$$

$$\begin{aligned}
 d(P, \tilde{C}) &\leq d((S, w), \tilde{C}) + \varepsilon \cdot \text{OPT} \\
 &\leq d((S, w), C^*) + \varepsilon \cdot \text{OPT} \\
 &\leq d(P, C^*) + 2\varepsilon \cdot \text{OPT} \\
 &= (1 + 2\varepsilon) \cdot \text{OPT}
 \end{aligned}$$

So it suffices to find a core set, then run approximation alg. on result.

Non-Streaming Coresets

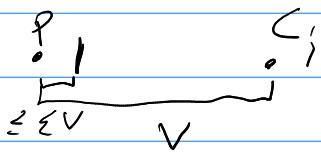
First, suppose we know good set of centers C . (e.g. from KP 99)

Want to get lots more centers so $d(P, S) \leq \epsilon \cdot d(P, C)$.
(then set weight = # mapping there)

1d:

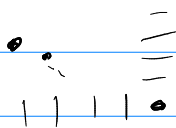
replace each c_i with $\{c_i \pm x \mid x = (1+\epsilon)^i\} \cup \{c_i\}$

$\cdot \times \times \times \times \times \times \times \times \times$



Previously contributed V
Now $\leq \epsilon V$.

2d: $\{c_i + (ax, bx) \mid x = 2^i, a \in [-\frac{1}{2}, \frac{1}{2}], b \in [-\frac{1}{2}, \frac{1}{2}]\}$



$\Rightarrow O\left(\frac{K}{\epsilon^2} \log \Delta\right)$ size for $O(\epsilon) \cdot OPT$.

(insertion)

Streaming Coresets

Merge-and-reduce (useful on $f(w_i)$)

Suppose $(S_1, w_1) = (k, \epsilon)$ coreset for P_1 ,
 $(S_2, w_2) = (k, \epsilon)$ coreset for P_2 .

Construct a coreset for $P_1 \cup P_2$.

Ans just run the algorithm, get (\bar{S}, \bar{w})
that is (k, ϵ') coreset.

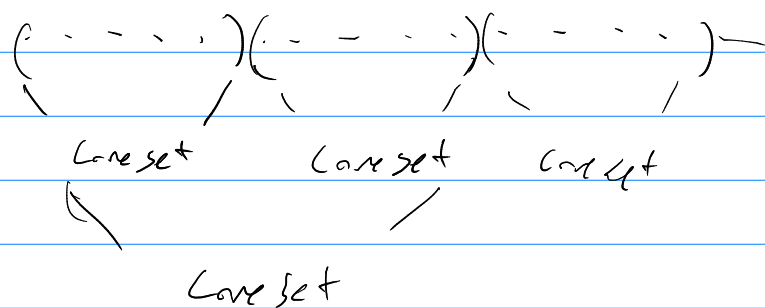
(S_1, w_1) moved P_1 by $\leq \epsilon \cdot \text{OPT}(P_1, k)$
so $(S_1 \cup S_2, w_1 \cup w_2)$ moved $P_1 \cup P_2$ by $\leq \epsilon \cdot \text{OPT}(P_1 \cup P_2, k)$

Hence $(S_1, w_1) \cup (S_2, w_2)$ is (k, ϵ) coreset for $P_1 \cup P_2$

But bigger. Run coreset alg. again to be smaller,
now

$$\begin{aligned} \text{dist}(\bar{S}, \bar{w}, P_1 \cup P_2) &\leq \epsilon \cdot \text{OPT}(P_1 \cup P_2, k) \\ &\quad + \text{dist}(\bar{S}, \bar{w}, (S_1 \cup S_2, w_1 \cup w_2)) \\ &\leq \epsilon \cdot \text{OPT}(P_1 \cup P_2, k) + \epsilon' \cdot \text{OPT}((S_1 \cup S_2, w_1 \cup w_2), k) \\ &\leq \epsilon \cdot \text{OPT}(P_1 \cup P_2, k) + \epsilon' (1 + \epsilon) \text{OPT}(P_1 \cup P_2, k) \\ &\leq (\epsilon + 2\epsilon') \text{OPT}(P_1 \cup P_2, k) \end{aligned}$$

Stream of Points



Store tree. Merge any two at same level to core set of size $O\left(\frac{k}{\epsilon^2} \log \Delta\right)$.

The approximation error has

$$\epsilon_i \leq 2\epsilon + \epsilon_{i-1} \leq 2^i \epsilon.$$

So final result is a $(k, 2\epsilon \cdot \log n)$ core set.
replace ϵ by $\frac{\epsilon}{2 \log n}$

$$\Rightarrow O\left(\frac{k}{\epsilon^2} \cdot \log^2 n \log \Delta\right) \text{ space per level}$$

$$O\left(\frac{k}{\epsilon^2} \log^3 n \log \Delta\right) \text{ total.}$$