

Lecture 22 Coresets

A coreset: "Turning big data into tiny data"

Idea: dramatically reduce the number of points while the cost of any solution is well approximated

I. A coreset for (Euclidean) k -means clustering

Let $D = D_{\ell_2^2}$. Let $D(A, C) := \sum_{x \in A} \min_{c \in C} \|x - c\|^2$.

Let $w: A \rightarrow \mathbb{R}^+$ be a weight function.

Let $D(A, w, C) := \sum_{x \in A} w(x) \cdot \min_{c \in C} \|x - c\|^2$.

Def: Let $A \subseteq \mathbb{R}^d$ be a set of n points, $\varepsilon \in (0, 1)$, and $k \in \mathbb{N}^{\geq 1}$. A set $S \subseteq \mathbb{R}^d$ together with a weight function $w: S \rightarrow \mathbb{R}^+$ is called a (k, ε) -coreset if for all $C \subseteq \mathbb{R}^d$ with $|C| \leq k$,

$$|D(A, C) - D(S, w, C)| \leq \varepsilon \cdot D(A, C).$$

Basic idea: Given A , first compute a discretization S of the space, snap each point to its nearest neighbor in S , and then use the (weighted) set S to approximate the cost function.

II Some tools:

[Def:] An ε -ball-cover of unit sphere (in Euclidean spaces) is a set of points B , such that for any point p in the unit sphere, the distance between p to B is at most ε .

Lemma 1: Let U be the unit sphere in d -dimensional Euclidean space. Then for every $0 < \varepsilon < 1$ there exists an ε -ball-cover B of size $(1 + \frac{2}{\varepsilon})^d$, i.e. for every point $p \in U$, $\min_{b \in B} \|p - b\| \leq \varepsilon$.

Remark: ① The above ball-cover is existential, There are algorithms for constructing one using $\varepsilon^{-O(d)}$ points.
② If not unit sphere, say U is a ball of radius r , then ε -ball-cover is the set B with $\min_{b \in B} \|p - b\| \leq \varepsilon \cdot r$.

Thm 1: There exists a 6.357-approximation algorithm for k -means problem. [Ahmadian et. al. 2017]

Remark: Any constant approximation will work for constructing a coresset.

Lemma 2 (Generalized Triangle Inequality) Let a, b, c be points in Euclidean space. Then for any $\varepsilon \in (0, 1)$, we have

$$|\|a-c\|^2 - \|b-c\|^2| \leq \frac{1^2}{\varepsilon} \|a-b\|^2 + 2\varepsilon \cdot \|a-c\|^2$$

$$|D(a, C) - D(b, C)| \leq \frac{1^2}{\varepsilon} D(a, b) + 2\varepsilon \cdot D(a, C)$$

II. A (non-efficient) alg for constructing a (k, ε) -coreset for k -means clustering

Input A, k

1. first compute a 10-approximation C to the k -means problem, which can be done in polynomial time. (Eg. By Thm 1)

Let C_1, \dots, C_k be the corresponding k clusters with centers c_1, \dots, c_k .

2. for each $j=1, \dots, k$,

Let $\bar{F} = G_j$

Let B^j be the ball with radius

$$r_j = \frac{1}{n} \sum_{x \in F} \|x - c_j\|^2 \text{ centered at } c_j \text{ and}$$

let S^j be the $\frac{\varepsilon}{192}$ -ball-cover of B^j for $i=0, \dots, \log_{10} n$

$$\text{let } S_j = \bigcup_{i=0}^{\log_{10} n} S^i$$

3. For each $x \in A$, let $B(x) \in \bigcup_{j=1}^k S_j$ be the nearest point in the union of all ball-covers. Let $S = \bigcup_{x \in A} B(x)$.

For each $y \in S$, let $w(y)$ be the number of $x \in A$ such that $B(x) = y$.

4. Output (S, w)

IV Analysis:

Lemma 3: Let F be a set of \sqrt{n} points. Let B^i be the ball with radius $r_i := \sqrt{\frac{2^i}{n} \sum_{x \in F} \|x\|^2}$ centered at the origin and let S^i be the $\frac{\varepsilon}{3}$ -ball-cover of B^i . Denote by $S = \bigcup_{i=0}^{\log_{10} n} S^i$. Then $\sum_{x \in A} \min_{s \in S} \|x - s\|^2 \leq \varepsilon^2 \sum_{x \in F} \|x\|^2$.

Proof: $F_{\text{close}} = \{y \in F \mid \|y\|^2 \leq \frac{1}{n} \sum_{x \in F} \|x\|^2\}$

$$F_{\text{far}} = F \setminus F_{\text{close}}.$$

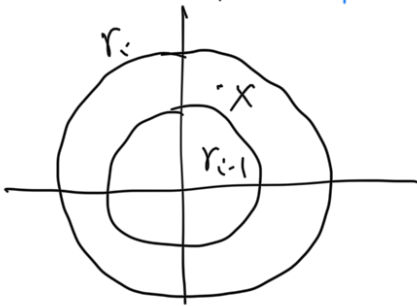
Since $|F_{\text{close}}| \leq n$, we have

$$\begin{aligned} \sum_{x \in F_{\text{close}}} \min_{s \in S^0} \|x - s\|^2 &\leq |F_{\text{close}}| \cdot \frac{1}{n} \sum_{x \in F} \|x\|^2 \cdot \frac{\varepsilon^2}{9} \\ &\leq \frac{\varepsilon^2}{9} \sum_{x \in F} \|x\|^2. \end{aligned}$$

For points in F_{far} , consider any point x in $B^i \setminus B^{i-1}$ for $i \in \{1, \dots, \log \log n\}$. We have

$$\min_{s \in S^i} \|x - s\|^2 \leq \frac{\varepsilon^2}{9} \cdot r_i^2 \leq \frac{4\varepsilon^2}{9} \cdot r_{i-1}^2 \leq \frac{4\varepsilon^2}{9} \cdot \|x\|^2$$

Summing up over all points, we have $\sum_{x \in F} \min_{s \in S} \|x - s\|^2$
 $\leq \frac{\varepsilon^2}{9} \cdot \sum_{x \in F} \|x\|^2 + \frac{4\varepsilon^2}{9} \sum_{x \in F} \|x\|^2 < \varepsilon^2 \sum_{x \in F} \|x\|^2$



the best $\text{poly}(k, \varepsilon^{-1})$.

Thm 1. The alg described above finds a coresnet for k -means consisting of $O(k \varepsilon^{-d} \log n)$ points, where d is the (constant) dimension.

Proof: For each $F \in \{C_1, \dots, C_k\}$, we have a total of $\log \log n$ balls with varying radii. For each such ball of radius r , we compute an $\frac{\varepsilon}{48}$ -ball-cover.

For any point $x \in A$, let $B(x)$ be the nearest point in the union of all ball-covers.

By Lemma 3, we have $\sum_{x \in A} \|x - B(x)\|^2 = \sum_{i=1}^k \sum_{x \in C_i} \|x - B(x)\|^2$

$$\leq \left(\frac{\varepsilon}{16}\right)^2 \sum_{i=1}^P \sum_{x \in C_i} \|x - c_i\|^2$$

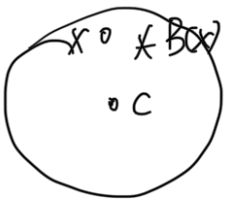
$$\leq \left(\frac{\varepsilon}{16}\right)^2 \cdot 10 \cdot D(A, C^*)$$

↑ i.e. OPT

Recall the definition of S, w from the algorithm.

Now consider an arbitrary set of centers C . We have

$$\begin{aligned} & |D(A, C) - D(S, w, C)| \\ &= \left| \sum_{x \in A} \min_{c \in C} \|x - c\|^2 - \sum_{x \in A} \min_{c \in C} \|B(x) - c\|^2 \right| \\ &\leq \sum_{x \in A} \left| \min_{c \in C} \|x - c\|^2 - \min_{c \in C} \|B(x) - c\|^2 \right| \\ &\leq \frac{12}{\varepsilon} \sum_{x \in A} \|x - B(x)\|^2 + 2\varepsilon \sum_{x \in A} \min_{c \in C} \|x - c\|^2 \\ &\leq \frac{12}{\varepsilon} \left(\frac{\varepsilon}{16}\right)^2 \cdot 10 \cdot D(A, C^*) + 2\varepsilon \cdot \sum_{x \in A} \min_{c \in C} \|x - c\|^2 \\ &< 2 \cdot \varepsilon \cdot D(A, C^*) + 2\varepsilon \cdot \sum_{x \in A} \min_{c \in C} \|x - c\|^2 \\ &\leq 4\varepsilon \cdot \sum_{x \in A} \min_{c \in C} \|x - c\|^2 \end{aligned}$$



Rescaling ε by a factor of $1/4$ finishes the proof. \square

Space: $O(k \cdot \log n)$ ball-covers, each with $O(\varepsilon^{-d})$ points.

V: Constructing coresets in the streaming (i.e. points arrive sequentially; limited storage)

Merge & Reduce technique.

① Composability

Lemma: Let $A_1, A_2 \subseteq \mathbb{R}^d$ be two ^{disjoint} point sets.

Assume S_1 with $w_1: S_1 \rightarrow \mathbb{R}$ and S_2 with $w_2: S_2 \rightarrow \mathbb{R}$ are (k, ε) -coresets for A_1, A_2 , respectively. Then $S_1 \cup S_2$ with $w_1 + w_2: S_1 \cup S_2 \rightarrow \mathbb{R}$ is a (k, ε) -coreset for $A_1 \cup A_2$.

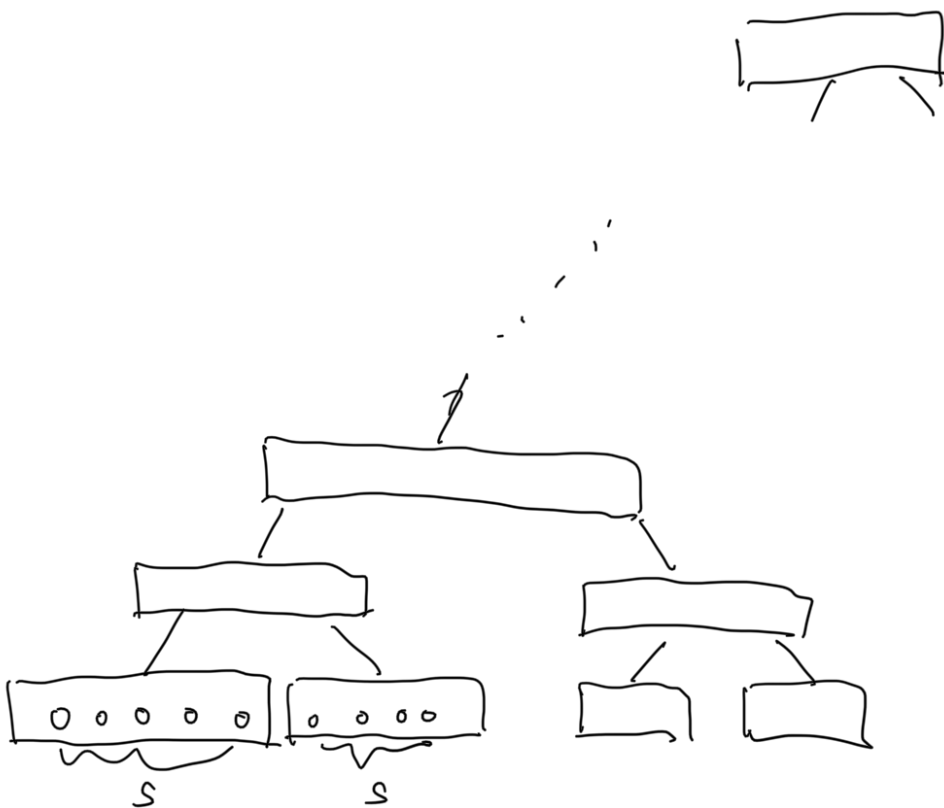
② Idea:

Assume n points in the sequence are partitioned into $\frac{n}{s}$ batches, each of size s .

These batches form the leaves of a binary tree of height $h \leq \log(n/s)$

For each batch, we compute a coresets
 Whenever we have computed a coresets for two children of a node, we aggregate them by recomputing a coresets of the union of the children and storing it in the parent node. The children can be then deleted at this point.

We only need store ≤ 2 coresets at each level.



Note: the final output (i.e. coresets at the root node) has approximation $(1 + \epsilon')^{\log n}$, where ϵ' is the parameter for each leaf. So we set $\epsilon' = \frac{\epsilon}{2(\log n + 1)}$, then

$$(1+\varepsilon')^{\log n} < 1+\varepsilon.$$

③ The algorithm

Initialize(k, ε, s)

1. Initialize an empty coresets B and initialize an empty dynamic array A
2. Set $\varepsilon' = \varepsilon / (2(\log n + 1))$, $i_{\max} = 0$
3. Store k, s .

Update($x \in \mathbb{R}^d$)

1. Store x in B
2. If $|B| = s$ do
3. Set $i = 0$
4. While $A[i] \neq \emptyset$ do
5. Compute a (k, ε') -coreset for $(A[i] \cup B, k)$, store it in B and empty $A[i]$
6. Set $A[i] = \emptyset$, $i = i + 1$,
 update $i_{\max} = \max\{i, i_{\max}\}$
7. EndWhile
8. Set $A[i] = B$ and empty B
9. EndIf

After any update, we can call the following `query()` to get the coresets for the points read so far.

`query()`

1. Set $T = B$
2. For $i = 0, \dots, i_{\max}$, do
3. Set $T = T \cup A[i]$
4. Compute a (k, ϵ') -coreset S for (T, k)
5. Return S

Thm: The above alg :

Space $O(k \epsilon^{-d} (\log n)^{d+2})$

(k, ϵ) -coreset of size $O(k (\log n)^{d+1} \cdot \epsilon^{-d})$.