

Coresets

---

## Coreset

(informally): a small

subset of the input

that "preserves" useful  
properties of the input

ideally  
very small  
compared to  
 $n$

## Example

• Minimum Enclosing Ball problem

Given:  $P = \{p_1, \dots, p_n\} \subseteq \mathbb{R}^d$

Find:  $x \in \mathbb{R}^d$  minimizing

$$\max_{p \in P} \{d(p, x)\}$$

Coreset:  $Q \subseteq P$  such that for all

$$y \in \mathbb{R}^d, \max_{q \in Q} \{d(q, y)\} \approx \max_{p \in P} \{d(p, y)\}.$$

## Coreset

(formally): given a cost function

$S$  is an  $\varepsilon$ -coreset for  $P$  if

- $(1-\varepsilon) \text{cost}(P) \leq \text{cost}(S) \leq \text{cost}(P)$
- $(1-\varepsilon) \text{cost}(P) \leq \text{cost}(S) \leq (1+\varepsilon) \text{cost}(P)$
- $\text{cost}(S) \leq \text{cost}(P) \leq (1+\varepsilon) \text{cost}(S)$

depends  
on the  
problem...

## Example

• Minimum Enclosing Ball problem

Given:  $P = \{p_1, \dots, p_n\} \subseteq \mathbb{R}^d$

Find:  $x \in \mathbb{R}^d$  minimizing

$$\text{cost}(P, x) = \max_{p \in P} \{d(p, x)\}$$

Coreset:  $Q \subseteq P$  such that for all  
 $y \in \mathbb{R}^d$ ,  $(1-\varepsilon)\text{cost}(P, y) \leq \text{cost}(Q, y) \leq \text{cost}(P, y)$

# $\epsilon$ -coreset for MEB

Alg:

$$S_1 \leftarrow \{ \text{an arbitrary } p \in P \}$$

for  $i = 1 + T \leftarrow \begin{matrix} \text{determine} \\ \text{later} \end{matrix}$

$$c_i \leftarrow \text{MEB center of } S_i$$

$$p_i \leftarrow \arg \max_{p \in P} d(c_i, p)$$

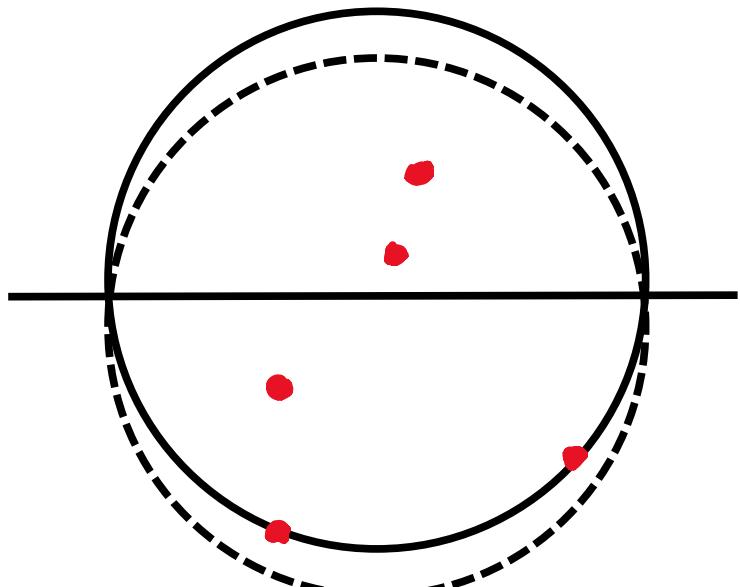
$$S_{i+1} \leftarrow S_i \cup \{p_i\}$$

return  $S_{T+1}$

# $\epsilon$ -coreset for MEB

Analysis:

**Lemma:** Let  $B$  be the MEB for  $P$  with center  $C$  & radius  $R$ . Any halfspace containing  $C$  contains  $p \in P$  with  $d(p, C) = R$ .



**Pf:** Suppose not. Then can shift the ball perpendicular to halfplane  $\Rightarrow$  no points on boundary.

We can then shrink ball  
 $\Rightarrow$  not minimum.

## $\epsilon$ -coreset for MEB

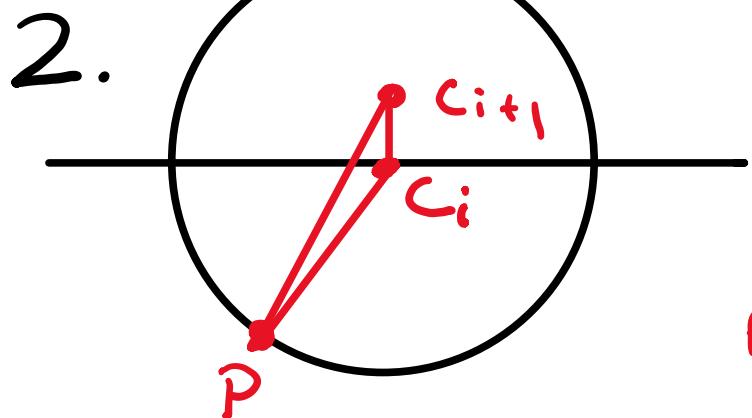
Analysis :  $R = \text{radius of MEB}(P)$        $\lambda_i = \frac{r_i}{R}$

$r_i = \text{radius of MEB}(S_i)$

Goal: identify  $T$  such that  $\lambda_T \geq 1 - \epsilon$

1.  $\exists q \in P$  s.t.  $d(q, c_i) \geq R$

$$r_{i+1} \geq d(q, c_{i+1}) \geq d(q, c_i) - d(c_i, c_{i+1}) \\ \geq R - d(c_i, c_{i+1})$$



Suppose  $d(c_i, c_{i+1}) > 0$  (else done)  
By Lemma  $\exists p, d(p, c_i) = r_i$

$$r_{i+1} \geq d(c_{i+1}, p) \geq \sqrt{r_i^2 + d(c_i, c_{i+1})^2}$$

## $\epsilon$ -coreset for MEB

Analysis :  $R = \text{radius of MEB}(P)$        $\lambda_i = \frac{r_i}{R}$   
 $r_i = \text{radius of MEB}(\mathcal{S}_i)$

Goal: identify  $T$  such that  $\lambda_T \geq 1 - \epsilon$

$$\Rightarrow \lambda_{i+1} \geq \frac{1}{R} \max \left\{ R - d(c_i, c_{i+1}), \sqrt{\lambda_i^2 R^2 + d(c_i, c_{i+1})^2} \right\}$$

minimized when  $R - d(c_i, c_{i+1}) = \sqrt{\lambda_i^2 R^2 + d(c_i, c_{i+1})^2}$

$$\Rightarrow d(c_i, c_{i+1}) = \frac{(1 - \lambda_i^2)R}{2}$$

$$\lambda_{i+1} \geq \frac{R - d(c_i, c_{i+1})}{R} = \frac{1 + \lambda_i^2}{2} \Rightarrow \lambda_i \geq 1 - \frac{1}{1 + i/2}$$

$\varepsilon$ -coreset for MEB

Analysis :  $R = \text{radius of MEB}(P)$        $\lambda_i = \frac{r_i}{R}$   
 $r_i = \text{radius of MEB}(S_i)$

Goal: identify  $T$  such that  $\lambda_T \geq 1 - \varepsilon$

$$\Rightarrow \lambda_T \geq 1 - \frac{1}{1 + T/2} \geq 1 - \varepsilon$$

Suffices to set  $T \geq \frac{2}{\varepsilon}$ .

# $\epsilon$ -coreset for MEB

Alg:

$$S_1 \leftarrow \{ \text{an arbitrary } p \in P \}$$

$$\text{for } i = 1 \text{ to } T = 2/\epsilon$$

$$c_i \leftarrow \text{MEB center of } S_i$$

$$p_i \leftarrow \arg \max_{p \in P} d(c_i, p)$$

$$S_{i+1} \leftarrow S_i \cup \{p_i\}$$

return  $S_{T+1}$

size of  
 $\epsilon$ -coreset is  
independent of  $n$

# Streaming Coresets

Useful properties:

- Reduce: If  $R$   $\epsilon$ -coreset for  $Q$   
 $Q$   $\delta$ -coreset for  $P$   
Then  $R$   $(\epsilon + \delta)$ -coreset for  $P$

Some version of this is always true

- (Strong) Merge: If  $P \cap P' = \emptyset$ ,  
 $Q$   $\epsilon$ -coreset for  $P$ ,  $Q'$   $\epsilon'$ -coreset for  $P'$   
 $Q \cup Q'$   $\max\{\epsilon, \epsilon'\}$ -coreset for  $P \cup P'$   
not always true. True for NEB

# Streaming Coresets

Reduce + Strong Merge Properties

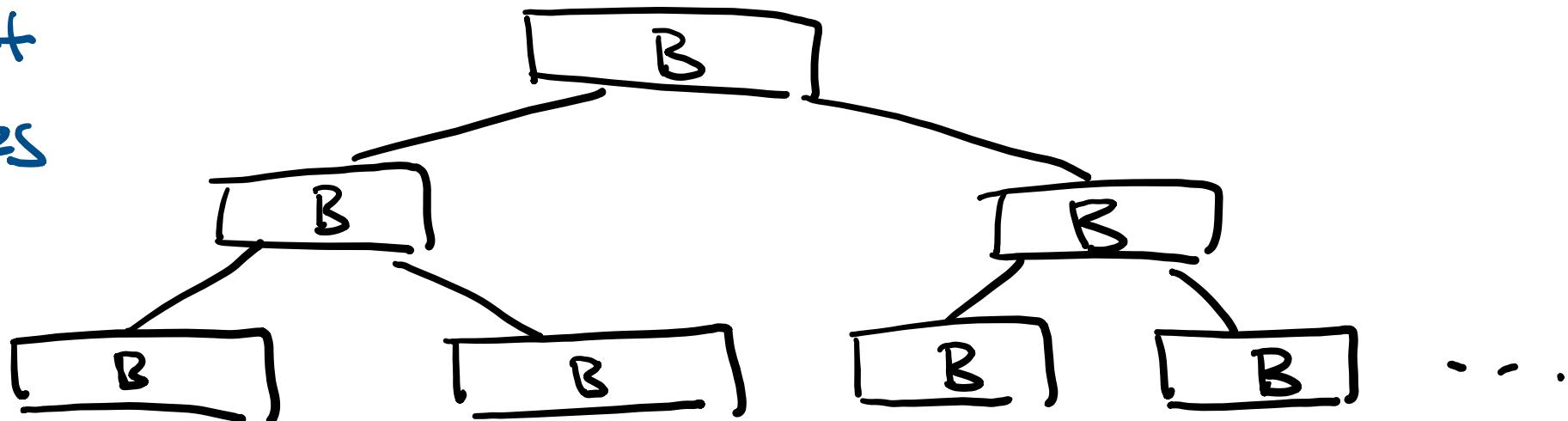
⇒ use Merge & Reduce Paradigm

Suppose  $f(\epsilon)$  space for  $\epsilon$ -coreset.

Set  $B = f(\epsilon/\log n)$

at each level:

$\epsilon/\log n$  -coreset  
of child nodes



# Streaming Coresets

Reduce + Strong Merge Properties

⇒ use Merge & Reduce Paradigm

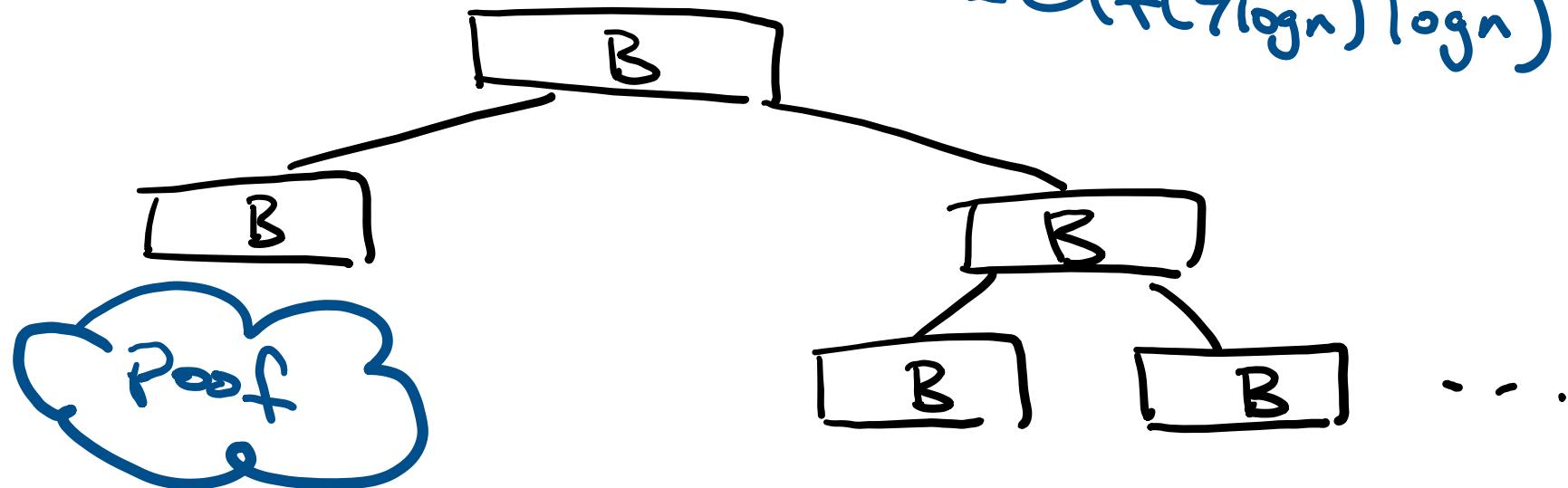
Suppose  $f(\varepsilon)$  space for  $\varepsilon$ -coreset.

$$\text{Set } B = f(\varepsilon / \log n)$$

at each level:

$\varepsilon / \log n$  -coreset  
of child nodes

keep at most  
 $2B$  points at  
each level



# Streaming Coresets

Reduce + Strong Merge Properties

⇒ use Merge & Reduce Paradigm

Suppose  $f(\varepsilon)$  space for  $\varepsilon$ -coreset.

Correctness?

by induction.

j-th level is  $j\frac{\varepsilon}{\log n}$  -coreset.

height =  $O(\log n)$  ⇒  $\varepsilon$ -coreset at top.

⇒  $O\left(\frac{(\log n)^2}{\varepsilon}\right)$  space for streaming MEB

$$\begin{aligned} \text{Space: } & O(B \log^n / B) \\ & = O(f(\varepsilon/\log n) \log n) \end{aligned}$$

# Clustering

k-center: Given  $P \subseteq \mathbb{R}^d$ ,  $k \in \mathbb{Z}$

find  $C = \{c_1, \dots, c_k\}$  minimizing

$$\text{cost}(P, C) = \max_{p \in P} \min_{c \in C} \{d(p, c)\}.$$

MEB:  $k$  is always 1.

k-median

$$\text{cost}(P, C) = \sum_{p \in P} \min_{c \in C} d(p, c)$$

k-means

$$\text{cost}(P, C) = \sum_{p \in P} \min_{c \in C} d(p, c)^2$$

NP-Hard  
if  $k$  is  
part  
of input

Coreset for  $k$ -median

$(k, \varepsilon)$ -coreset:  $Q \subseteq P$  such that

for all  $C \subseteq \mathbb{R}^d$ ,  $|C| = k$ ,

$$(1 - \varepsilon) \text{cost}(P, C) \leq \text{cost}(Q, C) \leq (1 + \varepsilon) \text{cost}(P, C)$$

will need idea of weighted  $k$ -median

$$\text{cost}(P, C) = \sum_{p \in P} w(p) \min_{c \in C} d(p, c)$$

Coreset for  $k$ -median

$(\alpha, \beta)$  - bicriterion approx.

Find  $A = \{a_1, \dots, a_m\}$      $m \leq \alpha k$

s.t.  $\text{cost}(P, A) \leq \beta \text{opt}(P, k)$

Will convert into  $(1+\varepsilon)$  -approx

using  $k$  "centers"

via weighted  $\varepsilon$ -coreset.

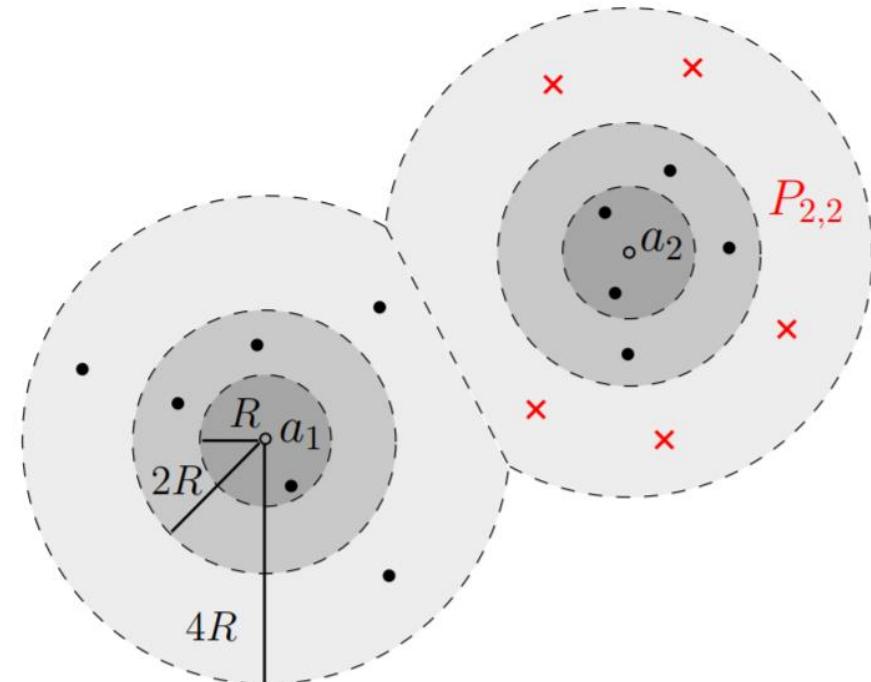
(w/ prob  $\geq 1-\delta$ )

# Coreset for k-median

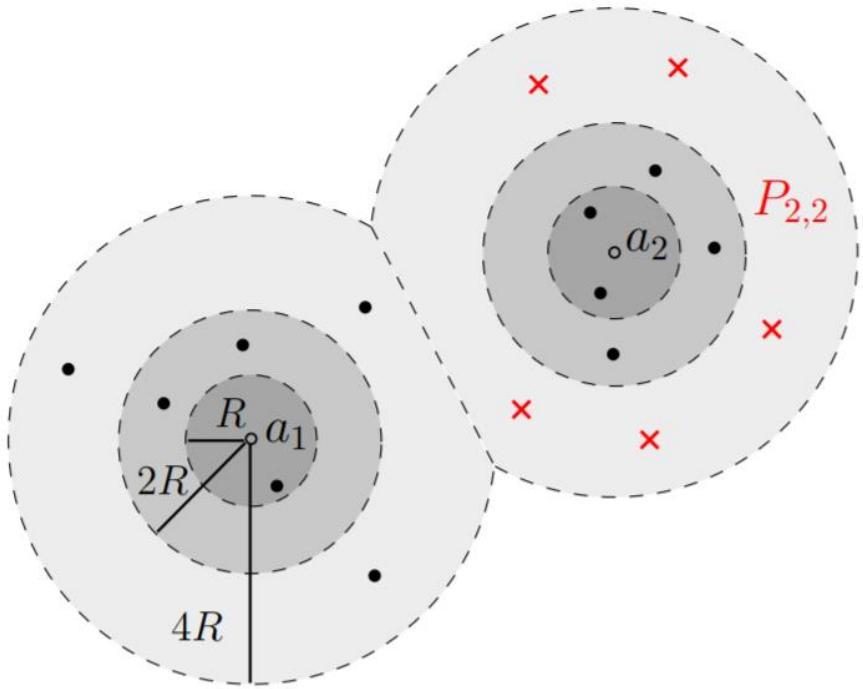
1. Let  $A$  be  $(\alpha, \beta)$  approx to k-median  
 $P_i$  : points in  $P$  "assigned" to  $a_i$
2. Split  $P_i$  into  $P_{i,j}$ ,  $j \in \{0 \dots \lg(\beta n)\}$

$P_{i,j} =$  points in  $P_i$  w/ distance to  $a_i$   
between  $2^{j-1}R \leq 2^j R$

$$R = \frac{\text{cost}(P, A)}{\beta n}$$



# Coreset for k-median



3. Sample  $s$  points from  $P_{ij}$

$$s = O\left(\frac{\beta^2}{\varepsilon^2} (k \log n + \log \frac{1}{\delta})\right)$$

↳ assign weight  $\frac{|P_{ij}|}{s}$

$$\rightarrow S_{ij} \subseteq P_{ij}.$$

Claim:  $S = \bigcup_{i,j} S_{ij}$  is  $\varepsilon$ -coreset  
w/ prob  $\geq 1 - \delta$

Coreset for k-median

$$S = O\left(\frac{\beta^2}{\xi^2} \left( k \log n + \log \frac{1}{\delta} \right)\right)$$

Key Lemma:  $U = \text{sample } \frac{\ln(2/\Delta)}{\xi^2}$  points from  $V$

Give each sample pt weight  $\frac{|U|}{|V|}$

$$|\text{cost}(V, C) - \text{cost}(U, C)| \leq \xi |V| \text{diam}(V) \text{ w/ prob } \geq 1 - \Delta.$$

Sampling  $S_{ij}$  from  $P_{ij} \Rightarrow$  need bound on  
 $|P_{ij}| \text{diam}(P_{ij})$

"easy" calculation:  $\sum |P_{ij}| \text{diam}(P_{ij}) \leq 6\beta \text{opt}(P, k)$

Coreset for k-median

$$S = O\left(\frac{\beta^2}{\varepsilon^2} (k \log n + \log \frac{1}{\delta})\right)$$

$$|\text{cost}(P_{ij}, C) - \text{cost}(S_{ij}, C)| \leq \xi |P_{ij}| \text{diam}(P_{ij}) \text{ w/ prob } \geq 1 - \Delta.$$

$$\sum |P_{ij}| \text{diam}(P_{ij}) \leq 6\beta \text{opt}(P, k)$$

$$\text{Set } \xi = \frac{\varepsilon}{6\beta}, \quad \Delta = \frac{n^{-2k} \delta}{2} \Rightarrow S = \frac{36\beta^2}{\varepsilon^2} \ln \frac{4n^{2k}}{\delta}$$

$$\begin{aligned} |\text{cost}(P, C) - \text{cost}(S, C)| &\leq \sum_{ij} |\text{cost}(P_{ij}, C) - \text{cost}(S_{ij}, C)| \\ &\leq \frac{\varepsilon}{6\beta} \sum_{ij} |P_{ij}| \text{diam}(P_{ij}) \leq \varepsilon \text{cost}(P, C) \end{aligned}$$

Coreset for k-median

$$s = O\left(\frac{\beta^2}{\varepsilon^2} \left( k \log n + \log \frac{1}{\delta} \right)\right)$$

now solve k-median directly on  $\varepsilon$ -coreset  
 $\Rightarrow (1+\varepsilon)$  -approximation.

Similar idea works for k-means.

!! cost functions for k-median & k-means  
satisfy strong merge  $\Rightarrow$  streaming !!

