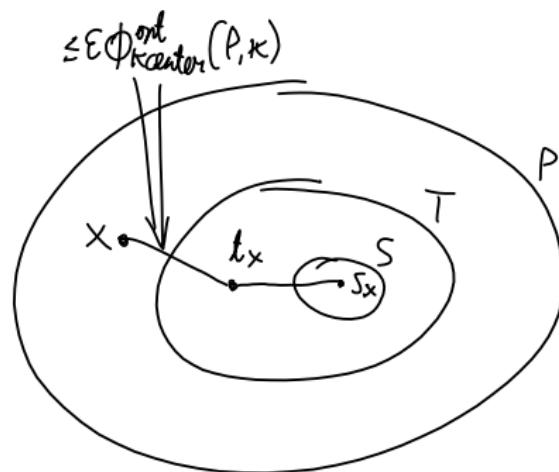


Coreset Technique

(Part 1 - Exercises)

Exercise

Let P be a set of N points in a metric space (M, d) , and let $T \subseteq P$ be a coresset of $|T| > k$ points such that for each $x \in P$ we have $d(x, T) \leq \epsilon \Phi_{k\text{center}}^{\text{opt}}(P, k)$, for some $\epsilon \in (0, 1)$. Let S be the set of k centers obtained by running the Farthest-First Traversal algorithm on T . Prove an upper bound to $\Phi_{k\text{center}}(P, S)$ as a function of ϵ and $\Phi_{k\text{center}}^{\text{opt}}(P, k)$.



$$\Phi_{k\text{center}}(P, S) = f(\epsilon) \cdot \Phi_{k\text{center}}^{\text{opt}}(P, k)$$

$$\forall x \in P: d(x, S) \leq f(\epsilon) \cdot \phi^{out}(P, \kappa)$$

t_x : punto vicino a x in T
 s_x : // // // a t_x in S

$$\begin{aligned} d(x, S) &\leq d(x, s_x) \leq d(x, t_x) + d(t_x, s_x) \\ &\leq \epsilon \phi^{out}(P, \kappa) + 2\phi^{out}(P, \kappa) \end{aligned}$$

$$d(x, S) \leq (2+\epsilon) \phi^{out}(P, \kappa) \Rightarrow (2+\epsilon) \text{ approx.}$$

Dimostriamo $\forall i \in \{1, \dots, k\}: d(x, z_i) \leq 2\phi^{out}(P, \kappa)$

$z = \{z_1, \dots, z_k\} \Rightarrow z_i$: punto vicinissimo a il. i di FFT

q : punto da it. extra di FFT $\Rightarrow q = z_{k+1} \Rightarrow S = S \cup \{q\} = \{z_1, \dots, z_{k+1}\}$

z e i cluster ottimali

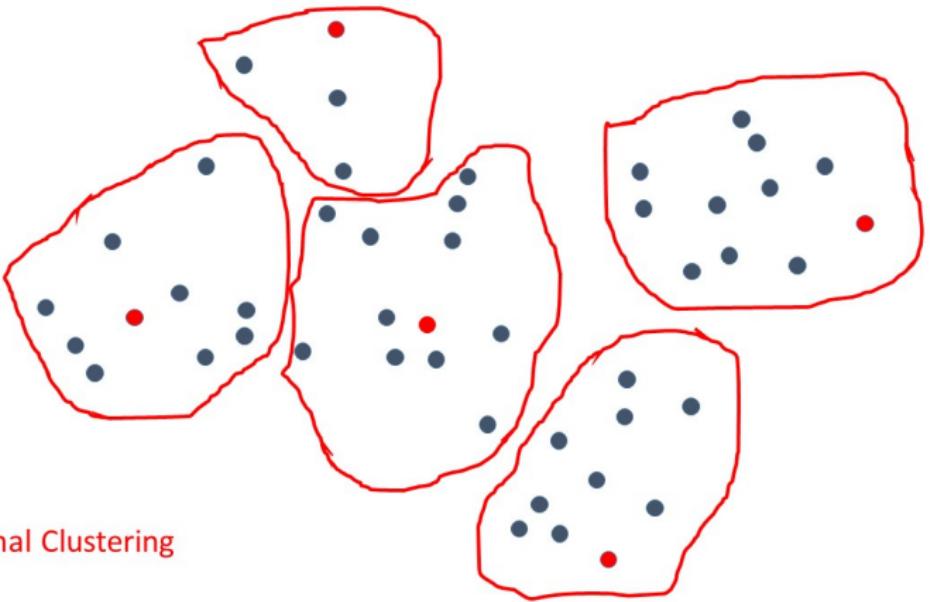
$\exists z_i, z_j, i < j \mid$ stereo cluster \hat{C}_e^* con centro \hat{x}_e^*

 $d(z_i, z_j) \leq d(z_i, \hat{c}_e^*) + d(z_j, \hat{c}_e^*) \leq \phi_{\text{out}}(P, k) + \phi_{\text{out}}(P, k) \leq 2\phi_{\text{out}}(P, k)$

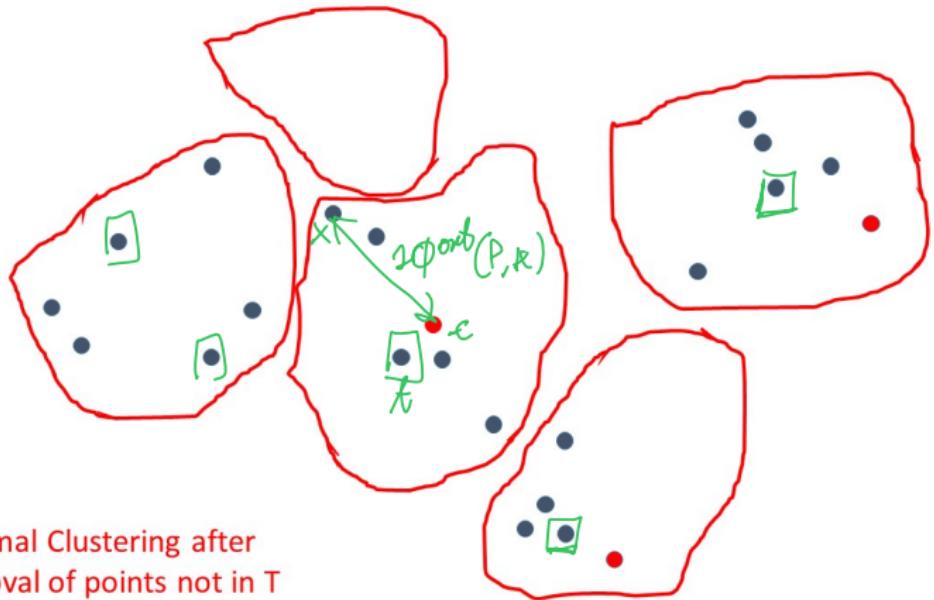
$$\forall x \in T, d(x, S) = d(x, \{z_1, \dots, z_k\}) \leq d(x, \{z_1, \dots, z_{j-1}\}) = \\ \leq d(z_j, \{z_1, \dots, z_{j-1}\}) \leq d(z_j, z_n) \leq 2\phi_{\text{out}}(P, k)$$

Exercise

Let P be a set of points in a metric space (M, d) , and let $T \subseteq P$.
For any $k < |T|, |P|$, show that $\Phi_{\text{kcenter}}^{\text{opt}}(T, k) \leq 2\Phi_{\text{kcenter}}^{\text{opt}}(P, k)$.
Is the bound tight?



Optimal Clustering



per ogni cluster, prendo punto
se ottengo $t \leq k$ punti, prendo altri $k - t$) } ottengo set T'

$$\Phi^{\text{ord}}(T, k) \leq \Phi(T, T) \leq 2\Phi^{\text{ord}}(P, k)$$

\downarrow
sol. generale prendendo
k punti

- 1) $C^* = \{C_1^*, \dots, C_R^*\}$ sol ottimale in P (centri $\hat{c}_1^*, \dots, \hat{c}_R^*$)
- 2) partiziono P usando C^*
- 3) solo altri punti non in T
- 4) $\forall C_i^*$ prendo 1 punto in $T \Rightarrow t_i$
- 5) $\hat{T} = \{\hat{t}_1, \dots, \hat{t}_R\}$, tutti punti $T \cap C_i^*$ assegnati a \hat{t}_i
- 6) $\Phi^*(T, k) \leq \Phi(T, \hat{T})$

$$\forall t \in T \quad d(t, \hat{T}) \leq d(t, \hat{t}_i) \leq d(t, \hat{c}_i^*) + d(\hat{c}_i^*, \hat{t}_i) \leq 2\Phi^{\text{ord}}(P, k)$$

(con $t \in C_i^*$)

Exercise

Let P be a set of N points in a metric space (M, d) , and let $\mathcal{C} = (C_1, C_2, \dots, C_k; c_1, c_2, \dots, c_k)$ be a k -clustering of P . Initially, each point $q \in P$ is represented by a pair $(\text{ID}(q), (q, c(q)))$, where $\text{ID}(q)$ is a distinct key in $[0, N - 1]$ and $c(q) \in \{c_1, \dots, c_k\}$ is the center of the cluster of q .

- ① Design a 2-round MapReduce algorithm that for each cluster center c_i determines the most distant point among those belonging to the cluster C_i (ties can be broken arbitrarily).
- ② Analyze the local and aggregate space required by your algorithm. Your algorithm must require $o(N)$ local space and $O(N)$ aggregate space.

ri può mandare tutti punti di un cluster al un reduce
no operazione su cluster, potremmo averne solo uno con tutti punti
 $M_L = \Theta(N)$

poniamo usare deterministic partitioning

ROUND 1

- Map: $(ID(q), (\alpha, \epsilon(q))) \rightarrow (ID(q) \bmod \sqrt{N}, (\alpha, \epsilon(q)))$
- Reduce: riunioni L_i : tutte parie con chiave $i =$ valore di tutte
parie con chiave i

• Sono c_i centri presenti in almeno una coppia in L_i
• $\forall c \in L_i$ calcolo punto + lontano da punti in L_i con centro c

• emetto $(c, q_c \wedge c \in L_i)$

ROUND 2:

- Map: niente
- Reduce: L_c valori da copie con chiave c
 - calcolo $\max_{q \in L_c} d(q, c) = d_{c, \text{MAX}}$
 - emetto $(c, d_{c, \text{MAX}})$

Exercise

Let P be a set of N bicolored points from a metric space, partitioned into k clusters C_1, C_2, \dots, C_k . Each point $x \in P$ is initially represented by the key-value pair $(\text{ID}_x, (x, i_x, \gamma_x))$, where ID_x is a distinct key in $[0, N - 1]$, i_x is the index of the cluster which x belongs to, and $\gamma_x \in \{0, 1\}$ is the color of x .

- ① Design a 2-round MapReduce algorithm that for each cluster C_i checks whether all points of C_i have the same color. The output of the algorithm must be the k pairs (i, b_i) , with $1 \leq i \leq k$, where $b_i = -1$ if C_i contains points of different colors, otherwise b_i is the color common to all points of C_i .
- ② Analyze the local and aggregate space required by your algorithm. Your algorithm must require $o(N)$ local space and $O(N)$ aggregate space.

