

Clustering with Same Cluster Queries

Problem → Let X be a subset of some Euclidean space \mathbb{R}^d .
 Let $G_x = \{G_1, \dots, G_k\}$ be a clustering of X . $\exists x_1, x_2 \in X$ belong to the same cluster is denoted by $x_1 \in G_i, x_2 \in G_j$. Define $n = |X|$ & $k = \text{no. of clusters}$.
 Clustering G_x is center-based if \exists set of centers $\mu_1, \dots, \mu_k \in \mathbb{R}^d$ such that the clustering corresponds to the Voronoi diagram over those center points. OR
 $\forall x \in X \& i \leq k, x \in G_i \Leftrightarrow i = \arg \min_j d(x, \mu_j)$. We also assume that centers μ^* corresponding to G^* are centers of mass of corresponding clusters.
 $\mu_i^* = \frac{1}{|G_i|} \sum_{x \in G_i} x$. So, oracle's clustering corresponds to the optimal clustering.

γ -Margin Property → Let X be set of points in metric-space M .
 Let $G_x = \{G_1, \dots, G_k\}$ be a center-based clustering of X induced by centers $\mu_1, \mu_2, \dots, \mu_k \in M$. G_x satisfies the γ -margin property iff
 $\forall i \in [k] \& \forall x \in G_i \& \forall y \in X/G_i$,
 $\gamma d(x, \mu_i) \leq d(y, \mu_i)$

Query function → For a clustering $C^* = \{G_1^*, G_2^*, \dots, G_k^*\}$, a C^* -oracle is a function f^* that answers queries according to the clustering.

$$\delta_{C^*}(x_1, x_2) = \begin{cases} \text{true} & \text{if } x_1 \in C^* \\ \text{false} & \text{else} \end{cases}$$

- This is called a semi-supervised active clustering framework (SSAC)
- ↳ An SSAC instance is determined by the tuple (X, d, C^*) .
1. An SSAC algorithm A is called a g -solver if for a family \mathcal{G} of instances, for every instance (X, d, C^*) , it can recover C^* by access to (X, d) & g ~~queries~~ queries to a C^* -oracle.
 2. It is a polynomial g -solver if its time complexity is polynomial in $n \cdot k$.
 3. \mathcal{G} admits an $O(g)$ query complexity if \exists an algo A such that it is a polynomial g -solver & instances $\in \mathcal{G}$.

Cluster assignment Query → Asks for the cluster index,

$$\delta_{C^*}(x) = ? \text{ iff } x \in C_i^*. \text{ This can be replaced with } k \text{ same-cluster queries.}$$

Lemma 1 → Let (X, d, C) be a clustering instance that satisfies the γ -margin property. Let μ be the set of centers of C . Let μ'_i such that

$$d(\mu_i, \mu'_i) \leq \gamma(C) \epsilon, \text{ where } \gamma(C) = \max_{x \in G} d(x, \mu_i)$$

Then if $\gamma \geq 1 + 2\epsilon \Rightarrow$

$$\forall x \in G, \forall y \in X \setminus G \Rightarrow d(x, \mu'_i) < d(y, \mu'_i)$$

(3)

Proof → Fix any $x \in G$ & $y \in G$. $d(x, \mu_i^*) \leq d(x, \mu_i) + d(\mu_i, \mu_i^*)$
 $\leq r(c_i)(1+\epsilon)$
 Similarly, $d(y, \mu_i^*) \geq d(y, \mu_i) - d(\mu_i, \mu_i^*) \geq (Y-\epsilon)r(c_i)$
 From both of them, $\Rightarrow d(x, \mu_i^*) < \frac{1+\epsilon}{Y-\epsilon} d(y, \mu_i^*) < d(y, \mu_i^*)$
 Hence proved.

Lemma 2 → Let Z_p, G, μ_p, μ_p' be defined as in the following algo, & $\epsilon = \frac{Y-1}{2}$. If $|Z_p| > \eta$ then the probability that $d(\mu_p, \mu_p') > r(c_p)\epsilon < \frac{\delta}{k}$

Theorem → Let (X, d, C) be a clustering instance, where C is center-based & satisfies the Y -margin property.
 Let $\mu_i = \frac{1}{|G_i|} \sum_{x \in G_i} x$. Assume $\delta \in (0, 1)$ & $Y > 1$. Then with probability $> 1 - \delta$, the following algorithm outputs C .

Algorithm → ^{Input} Clustering instance X , oracle Q , number of clusters k , parameters $\delta \in (0, 1)$

Output → A clustering C of set X ,

Initialize → $C = \{\}$, $S_1 = X$; $\eta = \beta \frac{\log k + \log(1/\delta)}{(Y-1)^2}$

~~Step~~ → Pseudocode → for $i=1$ to k do

Phase 1

$$l = k\eta + 1$$

$Z \sim U[S_i]$ // Draw l independent elements from S_i

For $1 \leq t \leq i$

$Z_t = \{x \in Z : Q(x) = t\}$ // Ask cluster assignment queries

$$P = \arg \max_t |Z_t|$$

$$\mu'_p = \frac{1}{|Z_p|} \sum_{x \in Z_p} x$$

Phase 2

// We know $\exists r_p$ such that $\forall x \in S_p, x \in G \Leftrightarrow d(x, \mu'_p) < r_p$
 // So r_p can be found by binary search.

$$\hat{S}_p = \text{Sorted}(\{S_i\}) \quad // \text{sort on basis of } d(x, \mu'_p)$$

$$r_p = \text{BinarySearch}(\hat{S}_p) \quad // \text{using up to } O(\log |S_p|) \text{ same cluster queries}$$

$$C' = \{x \in S_p : d(x, \mu'_p) \leq r_p\}$$

$$S_{A_1} = S_p \setminus C'$$

$$C = C \cup S_{A_1}$$

end

Proof → In first-phase, $l > k\eta$ cluster-assignment queries are made.

By pigeonhole-principle, $|Z_p| > \eta$. So, By lemma 2,

$d(\mu_p, \mu'_p) \leq \delta(C_p)$ with probability $> 1 - \frac{\delta}{k}$. By lemma 1,

⇒ $d(x, \mu'_p) < d(y, \mu'_p) \quad \forall x \in C_p \neq y \in C_p$. Hence radius r_p

found in phase 2 of algo is $r_p = \max_{x \in C_p} d(x, \mu'_p)$. Hence $C' = C_p$

So, with probability $\frac{\delta}{k}$, one iteration of algo

correctly finds all points in cluster C_p .

Let A_i denote the event that cluster $C_i \neq C'_i$, then

$$P(A_i) < \frac{\delta}{k}$$

② By union-bound

$$\Rightarrow P(A_1 \cup A_2 \cup \dots \cup A_k) < P(A_1) + P(A_2) + \dots + P(A_k) < k \left(\frac{\delta}{k} \right) = \delta$$

⇒ Probability of correct no-clustering $> [1 - \delta]$ hence proved.

5

Theorem → The algo makes $O(k \log n + k^2 \frac{\log k + \log(\frac{1}{\delta})}{(\gamma-1)^2})$ same cluster queries to the oracle \mathcal{Q} .
 It runs in $O(kn \log n + k^2 \frac{\log k + \log(\frac{1}{\delta})}{(\gamma-1)^2})$ time.

Proof → In each iteration → phase 2 has $O(\log n)$ queries & phase 1 has $O(kn)$ assignment queries during the ~~one~~ whole algo (by reusing results). ~~Some~~ assignment queries are same as k same-cluster queries.
 So we get $O(k \log n + k^2 n)$ queries in total.
 We sort S_i in phase 2, $\rightarrow O(n \log n)$ per iter.
 Phase 1 runs in $O(kn)$ time ~~per~~ per iter.
 So $O(kn \log n + k^2 n)$
Hence proved.

Proof of lemma 2 → Define a uniform distribution U over C_p . Then,
 μ_p & μ'_p are real & empirical means.
Statistics → Some inequality shows they are close enough.