

(1)

Paper - Exact recovery of Mangled Clusters

- * Relaxes the δ -margin property of Ashtiani et al to include ellipsoid clusterings & still uses only $O(\log n)$ queries of polynomial time to find the clustering.
- ↳ The problem is defined by a triple (X, k, δ) where $X \subset \mathbb{R}^d$ is a set of n points, $k \geq 2$ is the number of clusters & $\delta \in \mathbb{R} > 0$ is the margin. We assume \exists a latent clustering $C = \{G_1, G_2, \dots, G_k\}$ over input set X . We have access to an oracle $Scg(x, x')$ that answers 1 if x, x' belong to the same cluster & -1 otherwise. The goal is to recover C using as few queries as possible.
- * No algorithm can take less than n queries if clustering is arbitrary. So, we must assume some structure which is a margin condition.

Let $W \in \mathbb{R}^{d \times d}$ be a positive semidefinite matrix. Hence W induces the seminorm $\|x\|_W = \sqrt{x^T W x}$ & pseudo-metric

$$d_W(x, y) = \|x - y\|_W$$

Definition → Clustering Margin → A cluster C has margin $\delta > 0$ if \exists a PSD matrix $W = W(C)$ & a point $c \in \mathbb{R}^d$ such that $\forall y \notin C \quad \forall x \in C$, we have $d_W(y, c) > \sqrt{1+\gamma} d_W(x, c)$. If this holds for all clusters then the clustering C has margin δ .

Define $\Delta(\hat{C}, c) = \min_{S_k \in \mathcal{P}_k} \frac{1}{2n} \sum_{i=1}^k |G_i \Delta \hat{C}_{G_i}|$ where S_k is set of all permutations of $[k]$. The goal is to minimize $\Delta(\hat{C}, c)$, & use minimum queries to get $\Delta(\hat{C}, c) = 0$. The rank of a cluster

C , denoted by $\text{rank}(C)$, is the rank of the subspace spanned by its points. (2)

Theorem \rightarrow Any instance (X, k, Y) , whose latent clustering C has margin γ . Let $n = |X|$, $d \leq d$ be max. rank of a cluster in C & let $f(\gamma, d) = \max \left\{ 2^d, O\left(\frac{\gamma}{\gamma} \ln\left(\frac{\gamma}{\gamma}\right)\right) \right\}$.
 RECVR outputs C with probability $1 - \delta$ with high probability runs in time $O((k \ln n)(n + k^2 \ln k))$ using $O((k \ln n)(k^2 d^2 \ln k + f(\gamma, d)))$ same-cluster queries.

RECVR \rightarrow 1. sampling \rightarrow Draw points uniformly at random till, for some cluster C , we have a good probability ~~size~~ \mathcal{E} of size $\approx d^2$. Then, with bounds, any ellipsoid E containing \mathcal{E} contains atleast half of C .

2. Computing the Minimum Volume Enclosing Ellipsoid (MVEE) $\rightarrow E = E(\mathcal{S}_C)$. E contains atleast half of C & some more points from X/C . Now, we find & remove these points.

3. Tesselating the MVEE \rightarrow To recover $C \cap E$, partition E into $\left(\frac{d}{\gamma}\right)^d$ hyperrectangles, each being monochromatic. So, they are separated in $\left(\frac{d}{\gamma}\right)^d$ queries.

Theorem \rightarrow Suppose we are given a subset $\mathcal{S} \subseteq C$ where C is any unknown cluster. Then we can learn $C \cap E(\mathcal{S}_C)$ using $\max \left\{ 2^d, O\left(\frac{\gamma}{\gamma} \ln\left(\frac{\gamma}{\gamma}\right)\right) \right\}$ same-cluster queries where $\gamma = \text{rank}(C)$ & $E(\mathcal{S})$ is MVEE of \mathcal{S} .

The MVEE \rightarrow compute an ellipsoid close to $\text{Conv}(S_c)$.

A d-rounding of S is any ellipsoid satisfying

$\frac{1}{d} E \subseteq \text{Conv}(S_c) \subseteq E$. $E = g(S_c)$ is a d-rounding of S_c ; the d can be lowered to σ as well.

Monochromatic Tessellation \rightarrow Definition \rightarrow Monochromatic Subset \rightarrow

A set $B \subset \mathbb{R}^d$ is monochromatic w.r.t cluster C if it doesn't contain 2 points x, y with $x \in C$ & $y \notin C$.

If E can be divided into m monochromatic subsets, then we can find $\approx m$ CNE in m queries. Construction for $m = \left(\frac{d}{\gamma} \log d\right)^d$.

Lemma \rightarrow $\frac{\gamma}{C} < d_w(x, y)$ if $x \in C \neq y \notin C$

Let z be the point w.r.t. which margin of C holds.

$$d_w(y, z) > \sqrt{1+\gamma} d_w(x, z)$$

By triangle inequality, $d_w(y, x) \geq d_w(y, z) - d_w(x, z)$
 $> \sqrt{1+\gamma} - 1 (d_w(x, z))$

$$\text{So, for } \gamma \leq C^2 - 2C \Rightarrow 1 + \gamma \geq \left(1 + \frac{\gamma}{C}\right)^2$$

$$\Rightarrow d_w(y, x) > \sqrt{\left(\frac{1+\gamma}{C}\right)^2 - 1} = \frac{\gamma}{C} d_w(x, z)$$

\Rightarrow So, if $x, y \in X$ such that $x \in C \neq y \notin C$, then

$|x_i - y_i| \geq \frac{\gamma}{d}$ for some i . (This is scaled down by $d_w(x, z)$)

\Rightarrow For $\ell = 1 + \frac{\gamma}{d}$, Hyperrectangle with sides $[\beta_i, \beta_i \ell]$

is monochromatic.

(4)

Consider only the positive orthant of the ellipsoid. Let the semi-axes of E be the canonical basis of \mathbb{R}^d & its center μ be the origin. Let l_i be the length of i th semi-axis of E . We cover it with $\log_p\left(\frac{l_i}{\rho}\right)$ intervals of length ~~increasing then~~ increasing geometrically with ρ .

$$T_\rho = \left\{ [\beta_0, \beta_1], [\beta_1, \beta_2], \dots, [\beta_{i-1}, \beta_i] \right\} \text{ where } \beta_0 > 0, \rho \geq 1$$

Definition \rightarrow Let R_+ be the positive orthant of \mathbb{R}^d . The ~~tessellation~~ tessellation R of $E \cap R_+$ is the set of $(d+1)^d$ hyporectangles expressed in canonical basis $\{\mathbf{u}_1, \dots, \mathbf{u}_d\}$ of E : $R = T_1 \times \dots \times T_d$. If $\beta_0 \approx \frac{1}{d} l_i$, then the point $(\beta_1, \beta_2, \dots, \beta_d)$ lies well inside $\text{conv}(S_E)$. By setting β_0, ρ, b appropriately, R satisfies the following

- ① $|R| \leq \max \left\{ 1, O\left(\frac{d}{\rho} \ln \frac{d}{\rho}\right)^d \right\}$
- ② $E \cap R_+ \subseteq \bigcup_{R \in R} R$
- ③ For every $R \in R$, the set $R \cap E$ is monochromatic w.r.t C

Algorithm - TessellationLearn(X, S_E, Y)

- 1 compute $E \leftarrow E_j(S_E)$ or any other surrounding of S_E
- 2 compute $E_X \leftarrow X \cap E$
- 3 compute β_0, ρ, b as a function of γ, δ
- 4 for every $y \in E_X$ do
 - map y to $R(y)$
- 5 $x_c \leftarrow$ any point in E
- 6 while there is some unlabeled R do

(5)

$\text{label}(R) \leftarrow \text{SCB}(x_c, y)$, where y is any point s.t. $R(y) = R$
 return all y mapped to R such that $\text{label}(R) = +1$
 $R(y)$ represents the f^n that maps y to the hypercycle
 it belongs to.

Exact recovery of all clusters →

Algorithm $\text{RECUR}(X, k, \gamma, \epsilon)$

$$\hat{G}_1, \hat{G}_2, \dots, \hat{G}_k \leftarrow \emptyset$$

while $|X| > \epsilon n$ do

draw samples with replacement from $|X|$ until $|S_C| \geq b_0 l^2 \ln k$
 for some C .

$G_C \leftarrow \text{TessellationLearn}(X, S_C, \gamma)$

add G_C to the corresponding \hat{G}

$$X \leftarrow X \setminus G_C$$

return $\hat{C} = \{\hat{G}_1, \dots, \hat{G}_k\}$

Lemma 1 → The clustering \hat{C} returned by RECUR deterministically
 satisfy $\Delta(\hat{C}, C) \leq \epsilon$. & for $\epsilon < \frac{1}{n}$ $\Delta(\hat{C}, C) = 0$

Lemma 2 → $\text{RECUR}(X, k, \gamma, \epsilon)$ makes $O(k^3 \ln(\frac{1}{\epsilon}))$ same queries in
 expectation & for all fixed $a \geq 1$, $\text{RECUR}(X, k, \gamma, \epsilon)$ with
 probability at least $1 - n^{-a}$ makes $O(k^3 \ln k \ln n)$ queries &
 runs in time $O((k \ln n)(a + k^3 \ln k)) = \tilde{O}(kn + k^3)$

Proof → Lemma 3 → $\text{RECUR}(X, k, \gamma, \epsilon)$ makes atmost $8k \ln(\frac{1}{\epsilon})$
 rounds in expectation & for all fixed $a \geq 1$, with probability
 atleast $1 - n^{-a}$ performs at most $(8k + 6a\sqrt{k}) \ln n$ rounds

At each round, with probability $\frac{1}{2}$, a fraction of $\frac{1}{4k}$ of points are labeled & removed. So at each round size of X drops by $(1 - \frac{1}{8k})$ in expectation. Hence roughly $8k \ln(\frac{1}{\epsilon})$ rounds occur to drop $|X|$ below ϵn .

Query cost of RECUR → RECUR draws at most $bkd^e \ln k$ requires almost k queries $\Rightarrow O(k^2 \ln k)$ same cluster queries. Each cluster assignment Tessellation Learn requires $f(d, \gamma)$ queries ~~then~~ $= O(1)$ queries. So total expected queries = $\boxed{O(\cancel{k^3} k \ln k \ln n)}$

Running time of RECUR → Drawing samples needs $O(k^2 \ln k)$ for tessellation learn $\rightarrow O(|S_C|^{3.5} \ln |S_C|)$ for MVEE construction. Since $|S_C| = O(d^2 \ln k)$ by construction, $= \tilde{O}(1)$ computing $E_X = X \wedge E$ takes time $O(|X| \text{poly}(d)) = O(n)$. Copying $E_X \wedge X \wedge E$ takes $O(|X| \text{poly}(d))$.

For the classification, in time $O(|X \wedge E|)$ we can build a dictionary mapping every $R \in \mathcal{R}$ to set $R \wedge E_X$. Then each classification takes $O(1)$ time. To enumerate all positive R queries $O(|X \wedge E| \text{poly}(d))$ time. So combining with no. of rounds bound, with probability atleast $1 - n^{-1}$, runs in time $O(k \ln k \ln n + k^2 \ln k)$