
Exact Recovery of Mangled Clusters with Same-Cluster Queries (supplementary material)

Marco Bressan

Dept. of CS, Univ. of Milan, Italy
marco.bressan@unimi.it

Nicolò Cesa-Bianchi

DSRC & Dept. of CS, Univ. of Milan, Italy
nicolo.cesa-bianchi@unimi.it

Silvio Lattanzi

Google
silviol@google.com

Andrea Paudice

Dept. of CS, Univ. of Milan, Italy &
Istituto Italiano di Tecnologia, Italy
andrea.paudice@unimi.it

1 Ancillary results

1.1 VC-dimension of ellipsoids

For any PSD matrix M , we denote by $E_M = \{\mathbf{x} \in \mathbb{R}^d : d_M(\mathbf{x}, \boldsymbol{\mu}) \leq 1\}$ the $\boldsymbol{\mu}$ -centered ellipsoid with semiaxes of length $\lambda_1^{-1/2}, \dots, \lambda_d^{-1/2}$, where $\lambda_1, \dots, \lambda_d \geq 0$ are the eigenvalues of M . We recall the following classical VC-dimension bound (see, e.g., [3]).

Theorem 5. *The VC-dimension of the class $\mathcal{H} = \{E_M : M \in \mathbb{R}^d, M \succeq 0\}$ of (possibly degenerate) ellipsoids in \mathbb{R}^d is $\frac{d^2+3d}{2}$.*

1.2 Generalization error bounds

The next result is a simple adaptation of the classical VC bound for the realizable case (see, e.g., [5, Theorem 6.8]).

Theorem 6. *There exists a universal constant $c > 0$ such that for any family \mathcal{H} of measurable sets $E \subset \mathbb{R}^d$ of VC-dimension $d < \infty$, any probability distribution \mathcal{D} on \mathbb{R}^d , and any $\varepsilon, \delta \in (0, 1)$, if S is a sample of $m \geq c \frac{d \ln(1/\varepsilon) + \ln(1/\delta)}{\varepsilon}$ points drawn i.i.d. from \mathcal{D} , then for any $E^* \in \mathcal{H}$ we have:*

$$\mathcal{D}(E \triangle E^*) \leq \varepsilon \quad \text{and} \quad \mathcal{D}(E' \setminus E^*) \leq \varepsilon$$

with probability at least $1 - \delta$ with respect to the random draw of S , where E is any element of \mathcal{H} such that $E \cap S = E^ \cap S$, and E' is any element of \mathcal{H} such that $E^* \cap S \subseteq E' \cap S$.*

The first inequality is the classical PAC bound for the zero-one loss, which uses the fact that the VC dimension of $\{E \triangle E^* : E \in \mathcal{H}\}$ is the same as the VC dimension of \mathcal{H} . The second inequality follows immediately from the same proof by noting that, for any $E^* \in \mathcal{H}$ the VC dimension of $\{E \setminus E^* : E \in \mathcal{H}\}$ is not larger than the VC dimension of \mathcal{H} because, for any sample S and for any $F, G \in \mathcal{H}$, $(F \setminus E^*) \cap S \neq (G \setminus E^*) \cap S$ implies $F \cap S \neq G \cap S$.

1.3 Concentration bounds

We recall standard concentration bounds for non-positively correlated binary random variables, see [2]. Let X_1, \dots, X_n be binary random variables. We say that X_1, \dots, X_n are non-positively

correlated if for all $I \subseteq \{1, \dots, n\}$ we have:

$$\mathbb{P}(\forall i \in I : X_i = 0) \leq \prod_{i \in I} \mathbb{P}(X_i = 0) \quad \text{and} \quad \mathbb{P}(\forall i \in I : X_i = 1) \leq \prod_{i \in I} \mathbb{P}(X_i = 1) \quad (1)$$

Lemma 4 (Chernoff bounds). *Let X_1, \dots, X_n be non-positively correlated binary random variables. Let $a_1, \dots, a_n \in [0, 1]$ and $X = \sum_{i=1}^n a_i X_i$. Then, for any $\varepsilon > 0$, we have:*

$$\mathbb{P}(X < (1 - \varepsilon)\mathbb{E}[X]) < e^{-\frac{\varepsilon^2}{2}\mathbb{E}[X]} \quad (2)$$

$$\mathbb{P}(X > (1 + \varepsilon)\mathbb{E}[X]) < e^{-\frac{\varepsilon^2}{2+\varepsilon}\mathbb{E}[X]} \quad (3)$$

1.4 Yao's minimax principle

We recall Yao's minimax principle for Monte Carlo algorithms. Let \mathcal{A} be a finite family of deterministic algorithms and \mathcal{I} a finite family of problem instances. Fix any two distributions \mathbf{p} over \mathcal{I} and \mathbf{q} over \mathcal{A} , and any $\delta \in [0, 1/2]$. Let $\min_{A \in \mathcal{A}} \mathbb{E}_{I \sim \mathbf{p}}[C_\delta(I, A)]$ be the minimum, over every algorithm A that fails with probability at most δ over the input distribution \mathbf{p} , of the expected cost of A over the input distribution itself. Similarly, let $\max_{I \in \mathcal{I}} \mathbb{E}_{A \sim \mathbf{q}}[C_\delta(I, A)]$ be the expected cost of the randomized algorithm defined by \mathbf{q} under its worst input from \mathcal{I} , assuming it fails with probability at most δ . Then (see [4], Proposition 2.6):

$$\max_{I \in \mathcal{I}} \mathbb{E}_{\mathbf{q}}[C_\delta(I, A)] \geq \frac{1}{2} \min_{A \in \mathcal{A}} \mathbb{E}_{\mathbf{p}}[C_{2\delta}(I, A)] \quad (4)$$

2 Supplementary material for Section 5

2.1 Monochromatic Tessellation

We give a formal version of the claim about the monochromatic tessellation of Section 5:

Theorem 7. *Suppose we are given an ellipsoid E such that $\frac{1}{d\Phi}E \subset \text{conv}(S_C) \subset E$ for some stretch factor $\Phi > 0$. Then for a suitable choice of β_i, ρ, b , the tessellation \mathcal{R} of the positive orthant of E (Definition 3) satisfies:*

- (1) $|\mathcal{R}| \leq \max\{1, O\left(\frac{d\Phi}{\gamma} \ln \frac{d\Phi}{\gamma}\right)^d\}$
- (2) $E \cap \mathbb{R}_+^d \subseteq \cup_{R \in \mathcal{R}} R$
- (3) for every $R \in \mathcal{R}$, the set $R \cap E$ is monochromatic

In order to prove Theorem 7, we define the tessellation and prove properties (1-3) for $\gamma \leq 1/2$. For $\gamma > 1/2$ the tessellation is defined as for $\gamma = 1/2$, and one can check all properties still hold. In the proof we use a constant $c = \sqrt{5}$ and assume $\gamma < c^2 - 2c$, which is satisfied since $c^2 - 2c = 5 - 2\sqrt{5} > 1/2$.

First of all, we define the intervals T_i . The base i -th coordinate is:

$$\beta_i = \frac{\gamma}{c\sqrt{2d}} \frac{L_i}{\Phi d} \quad (5)$$

Note that, for all i ,

$$\frac{L_i}{\beta_i} = \frac{\Phi c d \sqrt{2d}}{\gamma} \quad (6)$$

Define:

$$\alpha = \frac{\gamma}{c\sqrt{2}\Phi d} \quad (7)$$

and let:

$$b = \max\left(0, \left\lceil \log_{1+\alpha}\left(\frac{c\Phi d \sqrt{2d}}{\gamma}\right) \right\rceil\right) \quad (8)$$

(The parameter ρ of the informal description of Section 5 is exactly $1 + \alpha$). Finally, define the interval set along the i -th axis as:

$$T_i = \begin{cases} \{[0, \beta_i]\} & \text{if } b = 0 \\ \{[0, \beta_i], (\beta_i, \beta_i(1 + \alpha)], \dots, (\beta_i(1 + \alpha)^{b-1}, \beta_i(1 + \alpha)^b)\} & \text{if } b \geq 1 \end{cases} \quad (9)$$

Proof of (1). By construction, $|T_i| = b + 1$. Thus, $|\mathcal{R}| = \prod_{i \in [d]} |T_i| = (b + 1)^d$. Thus, if $b = 0$ then $|\mathcal{R}| = 1$, else by (8) and 6,

$$b = \left\lceil \frac{\ln\left(\frac{c\Phi d\sqrt{2d}}{\gamma}\right)}{\ln(1 + \alpha)} \right\rceil \quad (10)$$

$$\leq \left\lceil \frac{2}{\alpha} \ln\left(\frac{c\Phi d\sqrt{2d}}{\gamma}\right) \right\rceil \quad \text{since } \ln(1 + \alpha) \geq \alpha/2 \text{ as } \alpha \leq 1 \quad (11)$$

$$= \left\lceil \frac{2\sqrt{2}c\Phi d}{\gamma} \ln \frac{c\Phi d\sqrt{2d}}{\gamma} \right\rceil \quad \text{definition of } \alpha \quad (12)$$

$$= O\left(\frac{d\Phi}{\gamma} \ln \frac{d\Phi}{\gamma}\right) \quad \text{since } d\Phi \geq 1, \gamma \leq 1/2 \quad (13)$$

in which case $|\mathcal{R}| = O\left(\frac{d\Phi}{\gamma} \ln \frac{d\Phi}{\gamma}\right)^d$. Taking the maximum over the two cases proves the claim.

Proof of (2). We show for any $\mathbf{x} \in E \cap \mathbb{R}_+^d$ there exists $R \in \mathcal{R}$ containing \mathbf{x} . Clearly, if $\mathbf{x} \in E \cap \mathbb{R}_+^d$, then $\langle \mathbf{x}, \mathbf{u}_i \rangle \in [0, L_i]$ for all $i \in [d]$. But T_i covers, along the i -th direction \mathbf{u}_i , the interval from 0 to

$$\beta_i(1 + \alpha)^b = \beta_i(1 + \alpha)^{\max(0, \lceil \log_{1+\alpha}(L_i/\beta_i) \rceil)} \geq \beta_i(1 + \alpha)^{\lceil \log_{1+\alpha}(L_i/\beta_i) \rceil} \geq L_i \quad (14)$$

Therefore some $R \in \mathcal{R}$ contains \mathbf{x} .

Proof of (3). Given any hyperrectangle $R \in \mathcal{R}$, we show that the existence of $\mathbf{x}, \mathbf{y} \in R \cap E$ with $\mathbf{x} \in C$ and $\mathbf{y} \notin C$ leads to a contradiction. For the sake of the analysis we conventionally set the origin at the center $\boldsymbol{\mu}$ of E , i.e. we assume $\boldsymbol{\mu} = \mathbf{0}$.

We define $E_{\text{in}} = \frac{1}{\Phi_d} E$ and let $M = U\Lambda U^\top$ be its PSD matrix, where $U = [\mathbf{u}_1, \dots, \mathbf{u}_d]$ and $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_d)$. Note that $\lambda_i = \frac{1}{\ell_i^2} = \frac{\Phi^2 d^2}{L_i^2}$ where $\ell_i = \frac{L_i}{\Phi_d}$ is the length of the i -th semiaxis of E_{in} . For any $R \in \mathcal{R}$, let R_i be the projection of R on \mathbf{u}_i (i.e. R_i is one of the intervals of T_i defined in (9)). Let $D = D(R) = \{i \in [d] : 0 \notin R_i\}$. We let U_D and U_{-D} be the matrices obtained by zeroing out the columns of U corresponding to the indices in $[d] \setminus D$ and D , respectively. Observe that if $\mathbf{x}, \mathbf{y} \in R \cap E$ then:

$$\langle \mathbf{x} - \mathbf{y}, \mathbf{u}_i \rangle^2 < \alpha^2 \langle \mathbf{x}, \mathbf{u}_i \rangle^2 \quad \forall i \in D \quad (15)$$

$$\langle \mathbf{x} - \mathbf{y}, \mathbf{u}_i \rangle^2 \leq \beta_i^2 \quad \forall i \notin D \quad (16)$$

Now suppose C has margin at least γ for some $\gamma \in (0, c^2 - 2c]$, and suppose $\mathbf{x}, \mathbf{y} \in R \cap E$ with $\mathbf{x} \in C$ and $\mathbf{y} \notin C$. Through a set of ancillary lemmata proven below, this leads to the absurd:

$$\frac{\gamma^2}{c^2} < d_W(\mathbf{y}, \mathbf{x})^2 \quad \text{Lemma 5} \quad (17)$$

$$\leq d_M(\mathbf{y}, \mathbf{x})^2 \quad \text{Lemma 6} \quad (18)$$

$$< \alpha^2 d_M(\mathbf{x}, \boldsymbol{\mu})^2 + \frac{\gamma^2}{2c^2} \quad \text{Lemma 7} \quad (19)$$

$$\leq \frac{\gamma^2}{2c^2} + \frac{\gamma^2}{2c^2} \quad \text{Lemma 8} \quad (20)$$

In the rest of the proof we prove the four lemmata.

Lemma 5. $\frac{\gamma}{c} < d_W(\mathbf{y}, \mathbf{x})$.

Proof. Let \mathbf{z} be the point w.r.t. which the margin of C holds. By the margin assumption,

$$d_W(\mathbf{y}, \mathbf{z}) > \sqrt{1 + \gamma} \quad \text{and} \quad d_W(\mathbf{x}, \mathbf{z}) \leq 1 \quad (21)$$

By the triangle inequality then,

$$d_W(\mathbf{y}, \mathbf{x}) \geq d_W(\mathbf{y}, \mathbf{z}) - d_W(\mathbf{x}, \mathbf{z}) > \sqrt{1 + \gamma} - 1 \quad (22)$$

One can check that for $\gamma \leq c^2 - 2c$ we have $1 + \gamma \geq (1 + \frac{\gamma}{c})^2$. Therefore

$$d_W(\mathbf{y}, \mathbf{x}) > \sqrt{(1 + \gamma/c)^2} - 1 = \frac{\gamma}{c} \quad (23)$$

as desired. \square

Lemma 6. $d_W(\cdot) \leq d_M(\cdot)$.

Proof. By the assumptions of the theorem, $E_{\text{in}} \subseteq \text{conv}_{\mu}(C)$. Moreover, by the assumptions on $d_W(\cdot)$, the unit ball of $d_W(\cdot)$ contains $\text{conv}(C)$. Thus, the unit ball of $d_W(\cdot)$ contains the unit ball of $d_M(\cdot)$. This implies $W \preceq M$, thus $\|\cdot\|_W \leq \|\cdot\|_M$ and $d_W(\cdot) \leq d_M(\cdot)$. \square

Lemma 7. $d_M(\mathbf{y}, \mathbf{x})^2 < \alpha^2 d_M(\mathbf{x}, \boldsymbol{\mu})^2 + \frac{\gamma^2}{2c^2}$.

Proof. We decompose $d_M(\mathbf{y}, \mathbf{x})^2$ along the colspaces of U_D and U_{-D} :

$$d_M(\mathbf{y}, \mathbf{x})^2 = \|M^{1/2}(\mathbf{y} - \mathbf{x})\|_2^2 \quad (24)$$

$$= \|M^{1/2}(\mathbf{y} - \mathbf{x})\|_{U_D U_D^\top}^2 + \|M^{1/2}(\mathbf{y} - \mathbf{x})\|_{U_{-D} U_{-D}^\top}^2 \quad (25)$$

Next, we bound the two terms of (25). To this end, we need to show that for all $D \subseteq [d]$ and $\mathbf{v} \in \mathbb{R}^d$:

$$\|M^{1/2} \mathbf{v}\|_{U_D U_D^\top}^2 = \sum_{i \in D} \lambda_i \langle \mathbf{v}, \mathbf{u}_i \rangle^2 \quad (26)$$

Let indeed $J_D = \text{diag}(\mathbf{1}_D)$ be the selection matrix corresponding to the indices of D . Then $U_D = U J_D$, and so $U^\top U_D = U^\top U J_D = J_D$. This gives:

$$\|M^{1/2} \mathbf{v}\|_{U_D U_D^\top}^2 = \mathbf{v}^\top (U \Lambda^{1/2} U^\top) U_D U_D^\top (U \Lambda^{1/2} U^\top) \mathbf{v} \quad \text{definition of } M \text{ and } \|\cdot\|. \quad (27)$$

$$= \mathbf{v}^\top U \Lambda^{1/2} J_D J_D \Lambda^{1/2} U^\top \mathbf{v} \quad \text{since } U^\top U_D = J_D \quad (28)$$

$$= \mathbf{v}^\top U J_D \Lambda^{1/2} \Lambda^{1/2} J_D U^\top \mathbf{v} \quad \text{since } \Lambda, J_D \text{ are diagonal} \quad (29)$$

$$= \mathbf{v}^\top U_D \Lambda U_D^\top \mathbf{v} \quad \text{since } U J_D = U_D \quad (30)$$

$$= \|U_D^\top \mathbf{v}\|_\Lambda^2 \quad \text{by definition} \quad (31)$$

$$= \sum_{i \in D} \lambda_i \langle \mathbf{v}, \mathbf{u}_i \rangle^2 \quad (32)$$

Now we can bound the first term of (25):

$$\|M^{1/2}(\mathbf{y} - \mathbf{x})\|_{U_D U_D^\top}^2 = \sum_{i \in D} \lambda_i \langle \mathbf{y} - \mathbf{x}, \mathbf{u}_i \rangle^2 \quad \text{by (32)} \quad (33)$$

$$< \alpha^2 \sum_{i \in D} \lambda_i \langle \mathbf{x}, \mathbf{u}_i \rangle^2 \quad \text{by (15)} \quad (34)$$

$$= \alpha^2 \|M^{1/2} \mathbf{x}\|_{U_D U_D^\top}^2 \quad \text{by (32)} \quad (35)$$

$$\leq \alpha^2 \|M^{1/2} \mathbf{x}\|_{U U^\top}^2 \quad (36)$$

$$= \alpha^2 \|M^{1/2} \mathbf{x}\|_2^2 \quad \text{since } U U^\top = I \quad (37)$$

$$= \alpha^2 d_M^2(\mathbf{x}, \boldsymbol{\mu}) \quad \text{since } \boldsymbol{\mu} = \mathbf{0} \quad (38)$$

And for the second term of (25), we have:

$$\|M^{1/2}(\mathbf{y} - \mathbf{x})\|_{U_{-D}U_{-D}^\top}^2 = \sum_{i \notin D} \lambda_i \langle \mathbf{y} - \mathbf{x}, \mathbf{u}_i \rangle^2 \quad \text{by (32)} \quad (39)$$

$$\leq \sum_{i \notin D} \lambda_i \beta_i^2 \quad \text{by (16)} \quad (40)$$

$$= \sum_{i \notin D} \frac{\Phi^2 d^2}{L_i^2} \left(\frac{\gamma}{c\sqrt{2d}} \frac{L_i}{\Phi d} \right)^2 \quad \text{by definition of } \lambda_i \text{ and } \beta_i \quad (41)$$

$$= \sum_{i \notin D} \frac{\gamma^2}{2dc^2} \quad (42)$$

$$\leq \frac{\gamma^2}{2c^2} \quad (43)$$

Summing the bounds on the two terms shows that $d_M(\mathbf{y}, \mathbf{x})^2 < \alpha^2 d_M(\mathbf{x}, \boldsymbol{\mu})^2 + \frac{\gamma^2}{2c^2}$, as claimed. \square

Lemma 8. $\alpha^2 d_M(\mathbf{x}, \boldsymbol{\mu})^2 \leq \frac{\gamma^2}{2c^2}$.

Proof. By construction we have $\mathbf{x} \in E$ and $E = \Phi d \cdot E_{\text{in}}$. Therefore $\frac{1}{\Phi d} \mathbf{x} \in E_{\text{in}}$, that is:

$$1 \geq d_M\left(\frac{1}{\Phi d} \mathbf{x}, \boldsymbol{\mu}\right)^2 = \frac{1}{\Phi^2 d^2} d_M(\mathbf{x}, \boldsymbol{\mu})^2 \quad (44)$$

where we used the fact that $d_M(\cdot, \boldsymbol{\mu})^2 = \|\cdot\|_M^2$ since $\boldsymbol{\mu} = \mathbf{0}$. Rearranging terms, this proves that $d_M(\mathbf{x}, \boldsymbol{\mu})^2 \leq \Phi^2 d^2$. Multiplying by α^2 , we obtain:

$$\alpha^2 d_M(\mathbf{x}, \boldsymbol{\mu})^2 \leq \left(\frac{\gamma}{\sqrt{2c}\Phi d} \right)^2 \Phi^2 d^2 = \frac{\gamma^2}{2c^2} \quad (45)$$

as desired. \square

The proof of the theorem is complete.

2.2 Low-stretch separators and proof of Theorem 3

In this section we show how to compute the separator of Theorem 3. In fact, computing the separator is easy; the nontrivial part is Theorem 3 itself, that is, showing that such a separator always exists.

To compute the separator we first compute the MVEE $E_J = (M^*, \boldsymbol{\mu}^*)$ of S_C (see Section 5). We then solve the following semidefinite program:

$$\begin{aligned} & \max_{\alpha \in \mathbb{R}, \boldsymbol{\mu} \in \mathbb{R}^d, M \in \mathbb{R}^{d \times d}} \alpha \\ \text{s.t. } & M \succeq \alpha M^* \\ & \langle M, (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top \rangle \leq 1 \quad \forall \mathbf{x} \in S_C \\ & \langle M, (\mathbf{y} - \boldsymbol{\mu})(\mathbf{y} - \boldsymbol{\mu})^\top \rangle > 1 \quad \forall \mathbf{y} \in S_{\bar{C}} \end{aligned} \quad (46)$$

where, for any two symmetric matrices A and B , $\langle A, B \rangle = \text{tr}(AB)$ is the usual Frobenius inner product, implying $\langle M, (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top \rangle = d_M(\mathbf{x}, \boldsymbol{\mu})^2$. In words, the constraint $M \succeq \alpha M^*$ says that E must fit into E_J if we scale E_J by a factor $\Phi = 1/\sqrt{\alpha}$. The other constraints require E to contain all of S_C but none of the points of $S_{\bar{C}}$. The objective function thus minimizes the stretch Φ of E .

In the rest of this paragraph we prove Theorem 3.

Proof of Theorem 3 (sketch). To build the intuition, we first give a proof sketch where the involved quantities are simplified. The analysis is performed in the latent space \mathbb{R}^d with inner product $\langle \mathbf{u}, \mathbf{v} \rangle = \mathbf{u}^\top W \mathbf{v}$. Setting conventionally $\mathbf{z} = \mathbf{0}$, C then lies in the unit ball \mathcal{B}_0 and all points of $X \setminus C$ lie outside $\sqrt{1+\gamma} \mathcal{B}_0$. For simplicity we assume $\gamma \ll 1$ so that $\sqrt{1+\gamma} \simeq 1+\gamma$, but we can easily extend the result to any $\gamma > 0$. Now fix the subset $S_C \subseteq C$, and let $E_J = E_J(S_C)$ be the MVEE of S_C . Observe the following fact: \mathcal{B}_0 trivially satisfies (1), but in general violates (2); in contrast, E_J trivially satisfies (2), but in general violates (1). The key idea is thus to “compare” \mathcal{B}_0 and E_J and take, loosely speaking, the best of the two. To see how this works, suppose for instance E_J has small radius, say less than $\gamma/4$. In this case, $E = E_J$ yields the thesis. Indeed, since the center $\boldsymbol{\mu}^*$ of E_J is in \mathcal{B}_0 , then any point of E is within distance $1 + \gamma/4 \leq \sqrt{1+\gamma}$ of the center of \mathcal{B}_0 , and lies inside $\sqrt{1+\gamma} \mathcal{B}_0$. Thus E_J separates S_C from $X \setminus C$, satisfying (1). At the other extreme, suppose E_J is large, say with all its d semiaxes longer than $\gamma/4$. In this case, $E = \mathcal{B}_0$ yields the thesis: indeed, by hypothesis E fits entirely inside $4/\gamma E_J$, satisfying (2). Unfortunately, the general case is more complex, since E_J may be large along some axes and small along others. In this case, both \mathcal{B}_0 and E_J fail to satisfy the properties. This requires us to choose the axes and the center of E more carefully. We show how to do this with the help of Figure 1.

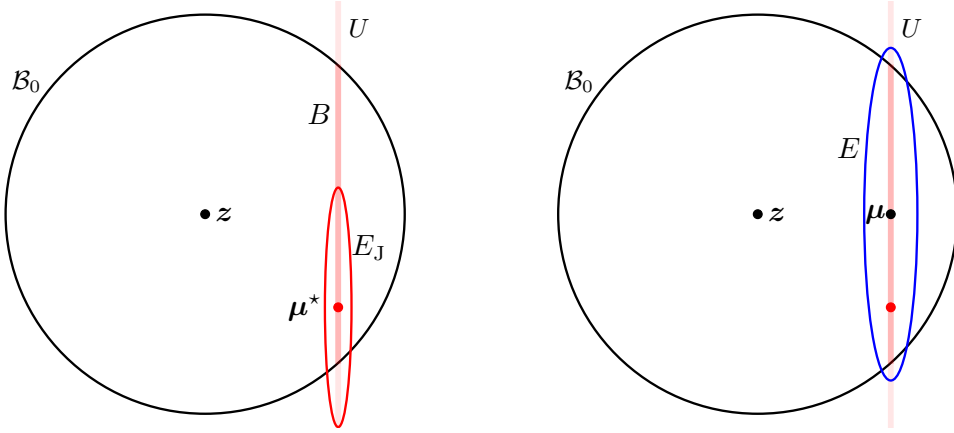


Figure 1: Left: the MVEE E_J of S_C and the affine subspace $U + \boldsymbol{\mu}^*$ (marked simply as U) spanned by its largest semiaxes. There is no guarantee that $E_J \subseteq \sqrt{1+\gamma} \mathcal{B}_0$. Right: the separator E , centered in the center $\boldsymbol{\mu}$ of B , with the largest semiaxis in U and the smallest one in U_\perp . We can guarantee that $S_C \subset E \subset \sqrt{1+\gamma} \mathcal{B}_0$.

Let $\{\mathbf{u}_1, \dots, \mathbf{u}_d\}$ be the orthonormal basis defined by the semiaxes of E_J and $\ell_1^*, \dots, \ell_d^*$ be the corresponding semiaxes lengths. We define a threshold $\varepsilon = \gamma^3/d^2$, and partition $\{\mathbf{u}_1, \dots, \mathbf{u}_d\}$ as $A_P = \{i : \ell_i^* > \varepsilon\}$ and $A_Q = \{i : \ell_i^* \leq \varepsilon\}$. Thus A_P contains the large semiaxes of E_J and A_Q the small ones. Let U, U_\perp be the subspaces spanned by $\{\mathbf{u}_i : i \in A_P\}$ and $\{\mathbf{u}_i : i \in A_Q\}$, respectively. Consider the subset $B = \mathcal{B}_0 \cap (\boldsymbol{\mu}^* + U)$. Note that B is a ball in at most d dimensions, since it is the intersection of a d -dimensional ball and an affine linear subspace of \mathbb{R}^d . Let $\boldsymbol{\mu}$ and ℓ be, respectively, the center and radius of B . We set the center of E at $\boldsymbol{\mu}$, and the lengths ℓ_i of its semiaxes as follows:

$$\ell_i = \begin{cases} \frac{\ell}{\sqrt{1+\gamma}} & \text{if } i \in A_P \\ \frac{\ell_i^*}{\sqrt{\varepsilon}} & \text{if } i \in A_Q \end{cases} \quad (47)$$

Loosely speaking, we are “copying” the semiaxes from either \mathcal{B}_0 or E_J depending on ℓ_i^* . In particular, the large semiaxes (in A_P) are set so to contain all of B and exceed it by a little, taking care of not intersecting $\sqrt{1+\gamma} \mathcal{B}_0$. Instead, the small semiaxes (in A_Q) are so small that we can safely set them to $1/\sqrt{\varepsilon}$ times those of E_J , so that we add some “slack” to include S_C without risking to intersect $\sqrt{1+\gamma} \mathcal{B}_0$. Now we are done, and our low-stretch separator is $(M, \boldsymbol{\mu})$ where $M = \sum_{i=1}^d \ell_i^{-2} \mathbf{u}_i \mathbf{u}_i^\top$. This the ellipsoid E that yields Theorem 3. In the next paragraph, we show how we can find efficiently all points in E that belong to C .

2.3 Proof of Theorem 3 (full).

We prove the theorem for $\gamma \leq 1/5$ and use the fact that whenever C has weak margin γ then it also has weak margin γ' for all $\gamma' > \gamma$. As announced, the analysis is carried out in the latent space \mathbb{R}^d equipped with the inner product $\langle \mathbf{u}, \mathbf{v} \rangle = \mathbf{u}^\top W \mathbf{v}$. All norms $\|\mathbf{u}\|$, distances $d(\mathbf{u}, \mathbf{v})$, and (cosine of) angles $\langle \mathbf{u}, \mathbf{v} \rangle / (\|\mathbf{u}\| \|\mathbf{v}\|)$ are computed according to this inner product unless otherwise specified. Let \mathcal{B}_0 be the unit ball centered at the origin, which we conventionally set at \mathbf{z} , the point in the convex hull of C according to which the margin is computed. Then, by assumption, $C \subset \mathcal{B}_0$, and $\mathbf{x} \notin \sqrt{1+\gamma} \mathcal{B}_0$ for all $\mathbf{x} \notin C$. For ease of notation, in this proof we denote the MVEE by E^* rather than E_J . Let then $(E^*, \boldsymbol{\mu}^*)$ be the MVEE of S_C ; note that $\boldsymbol{\mu}^* \in \text{conv}(S_C) \subseteq \mathcal{B}_0$. We let $\mathbf{u}_1, \dots, \mathbf{u}_d$ be the orthonormal eigenvector basis given by the axes of E^* and $\lambda_1^*, \dots, \lambda_d^*$ the corresponding eigenvalues. Note that if $\min_i \lambda_i^* \geq 5/\gamma^2$ then E^* has radius $\leq \gamma/\sqrt{5}$ and thus, since $\boldsymbol{\mu}^* \in \mathcal{B}_0$ and $\gamma \leq 1/5$, its distance from \mathcal{B}_0 is at most $1 + \gamma/\sqrt{5} = \sqrt{1 + 2\gamma/\sqrt{5} + \gamma^2/5} < \sqrt{1+\gamma}$. In this case we can simply set $E = E^*$ and the thesis is proven. Thus, from now on we assume $\min_i \lambda_i^* < 5/\gamma^2$.

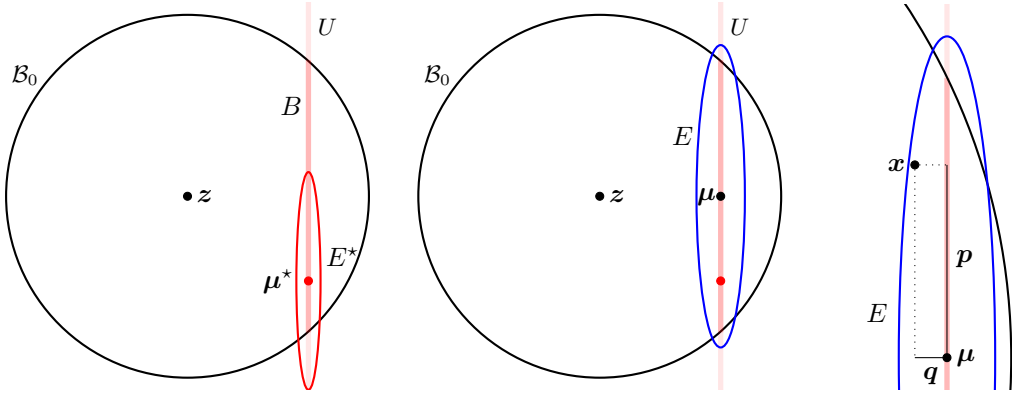


Figure 2: Left: the separating ball \mathcal{B}_0 of C , the MVEE E^* of S_C , and the affine subspace $U + \boldsymbol{\mu}^*$ spanned by its largest semi-axes. Middle: E is our separator centered in the center $\boldsymbol{\mu}$ of the ball $B = U \cap \mathcal{B}_0$. Right: a point $\mathbf{x} \in S_C$ with its projections onto U and U_\perp with respect to the origin, which we conventionally set at $\boldsymbol{\mu}$ (the center of E).

Now let:

$$\varepsilon = \frac{\gamma^3}{32d^2} \quad (48)$$

and partition (the indices of) the basis $\{\mathbf{u}_1, \dots, \mathbf{u}_d\}$ as follows:

$$A_P = \{i : \lambda_i^* < 1/\varepsilon^2\}, \quad A_Q = [d] \setminus A_P \quad (49)$$

Since $\min_i \lambda_i^* < 5/\gamma^2$ and $5/\gamma^2 \leq 1/\varepsilon^2$, then by construction the set A_P is not empty. We now define the ellipsoid E . Let U, U_\perp be the subspaces spanned by $\{\mathbf{u}_i : i \in A_P\}$ and $\{\mathbf{u}_i : i \in A_Q\}$ respectively, and let $B = \mathcal{B}_0 \cap (\boldsymbol{\mu}^* + U)$. Note that B is a ball, since it is the intersection of a ball and an affine linear subspace. Let $\boldsymbol{\mu}$ and ℓ be, respectively, the center and radius of B and define

$$\lambda_i = \begin{cases} (1 - \sqrt{5\gamma/4})\ell^{-2} & i \in A_P \\ \varepsilon\lambda_i^* & i \in A_Q \end{cases} \quad M = \sum_{i=1}^d \lambda_i \mathbf{u}_i \mathbf{u}_i^\top \quad (50)$$

Then our ellipsoidal separator is $E = \{\mathbf{x} \in \mathbb{R}^d : d_M(\mathbf{x}, \boldsymbol{\mu}) \leq 1\}$. See Figure 2 for a pictorial representation. We now prove that E satisfies: **(1)** $S_C \subset E$, **(2)** $E \subseteq \frac{64\sqrt{2}d^2}{\gamma^3} E^*(S_C)$, **(3)** $E \subset \sqrt{1+\gamma} \mathcal{B}_0$.

Proof of (1). Set the center $\boldsymbol{\mu}$ of E as the origin. For all $i \in [d]$ let $U_i = \mathbf{u}_i \mathbf{u}_i^\top$ and define the following matrices:

$$P_0 = \sum_{i \in A_P} U_i, \quad Q_0 = \sum_{i \in A_Q} U_i \quad (51)$$

$$P = \sum_{i \in A_P} \lambda_i U_i, \quad Q = \sum_{i \in A_Q} \lambda_i U_i \quad (52)$$

$$P_\star = \sum_{i \in A_P} \lambda_i^\star U_i, \quad Q_\star = \sum_{i \in A_Q} \lambda_i^\star U_i \quad (53)$$

We want to show that $d_M^2(\mathbf{x}, \boldsymbol{\mu}) \leq 1$ for all $\mathbf{x} \in S_C$. Note that $d_M(\mathbf{x}, \boldsymbol{\mu})^2$ equals (recall that $\boldsymbol{\mu} = \mathbf{0}$):

$$\mathbf{x}^\top P \mathbf{x} + \mathbf{x}^\top Q \mathbf{x} \quad (54)$$

Let us start with the second term of (54). By definition of Q_\star and since $\boldsymbol{\mu}^{\star\top} Q_\star = (\boldsymbol{\mu}^\star - \boldsymbol{\mu})^\top Q_\star = \mathbf{0}$ because $\boldsymbol{\mu}^\star - \boldsymbol{\mu} \in U$,

$$\mathbf{x}^\top Q \mathbf{x} = \varepsilon \mathbf{x}^\top Q_\star \mathbf{x} = \varepsilon (\mathbf{x} - \boldsymbol{\mu}^\star)^\top Q_\star (\mathbf{x} - \boldsymbol{\mu}^\star) \leq \varepsilon < \frac{\gamma}{4} \quad (55)$$

where the penultimate inequality follows from $\mathbf{x} \in E^\star$.

We turn to the first term of (54). If we let \mathbf{p}, \mathbf{q} be the projections of $\mathbf{x} - \boldsymbol{\mu} = \mathbf{x}$ onto U, U_\perp , so that

$$\|\mathbf{p}\|^2 = \mathbf{x}^\top P_0 \mathbf{x}, \quad \|\mathbf{q}\|^2 = \mathbf{x}^\top Q_0 \mathbf{x} \quad (56)$$

then by definition of the λ_i we have:

$$\mathbf{x}^\top P \mathbf{x} = \frac{1 - \sqrt{5\gamma/4}}{\ell^2} \|\mathbf{p}\|^2 \quad (57)$$

We can thus focus on bounding $\|\mathbf{p}\|$. Since B is a ball of radius ℓ , then $\|\mathbf{p}\| \leq \ell + d(\mathbf{p}, B)$, where $d(\mathbf{p}, B)$ is the distance of \mathbf{p} from its projection on B —see Figure 3, left.

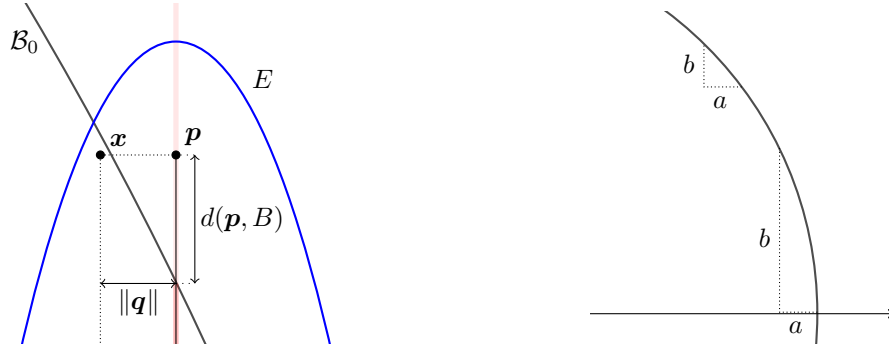


Figure 3: Left: a point $\mathbf{x} \in S_C \subset \mathcal{B}_0$ which lies in E as well. Right: for a fixed $a > 0$, the ratio b/a is maximized when the segment of length a lies on the line passing through the center of \mathcal{B}_0 , in which case $b/a = \frac{\sin \theta}{1 - \cos \theta}$ for some $\theta \in (0, \pi/2)$.

Now, since $\mathbf{x} \in \mathcal{B}_0$, the ratio $\frac{d(\mathbf{p}, B)}{\|\mathbf{q}\|}$ is maximized when $\ell \rightarrow 0$ (i.e., B has a vanishing radius), in which case $d(\mathbf{p}, B) \leq \sin \theta$ and $\|\mathbf{q}\| \geq 1 - \cos \theta$, where $\theta \in (0, \pi/2]$; see Figure 3 right. Then:

$$\frac{\|\mathbf{q}\|}{d(\mathbf{p}, B)} \geq \frac{1 - \cos \theta}{\sin \theta} = \tan \frac{\theta}{2} \geq \frac{\theta}{2} \geq \frac{\sin \theta}{2} \geq \frac{d(\mathbf{p}, B)}{2} \quad (58)$$

where we used the tangent half-angle formula and the Taylor expansion of $\tan \theta$. This yields $d(\mathbf{p}, B) \leq \sqrt{2\|\mathbf{q}\|}$. Thus:

$$\|\mathbf{p}\| \leq \ell + \sqrt{2\|\mathbf{q}\|} \quad (59)$$

But since $\lambda_i^* \geq 1/\varepsilon^2$ for all $i \in A_Q$:

$$\|\mathbf{q}\|^2 = \mathbf{x}^\top Q_0 \mathbf{x} \leq \varepsilon^2 \mathbf{x}^\top Q_* \mathbf{x} = \varepsilon^2 (\mathbf{x} - \boldsymbol{\mu}^*)^\top Q_* (\mathbf{x} - \boldsymbol{\mu}^*) \leq \varepsilon^2 \quad (60)$$

Therefore:

$$\mathbf{x}^\top P \mathbf{x} \leq \frac{1 - \sqrt{5\gamma/4}}{\ell^2} (\ell + \sqrt{2\varepsilon})^2 \leq (1 - \sqrt{5\gamma/4})(1 + \sqrt{2\varepsilon}/\ell)^2 \quad (61)$$

Next, we show that $\frac{\sqrt{2\varepsilon}}{\ell} \leq \frac{1}{2}\sqrt{5\gamma/4}$. First,

$$\sqrt{2\varepsilon} = \sqrt{2 \frac{\gamma^3}{32d^2}} = \frac{\gamma\sqrt{\gamma}}{4d} \quad (62)$$

We now temporarily set $\boldsymbol{\mu}^*$ as the origin. We want to show that the projection of $1/d E^*$ on U is contained in B . Now, the projection of an ellipsoid on the subspace spanned by a subset of its axes is a subset of the ellipsoid itself, and U is by definition spanned by a subset of the axes of E^* . Therefore the projection P of $1/d E^*$ on U satisfies $P \subseteq 1/d E^*$. Suppose then by contradiction that $P \not\subseteq B$. Since $B = U \cap \mathcal{B}_0$, this implies that $1/d E^* \not\subseteq \mathcal{B}_0$. But by John's theorem, $1/d E^* \subseteq \text{conv}(S_C)$, and therefore $\text{conv}(S_C) \not\subseteq \mathcal{B}_0$, which is absurd. Therefore $P \subseteq B$.

Let us get back to the proof, with $\boldsymbol{\mu}$ as the origin. On the one hand, the definitions of A_P and U imply that the largest semiaxis of E^* of length $\ell^* = 1/\sqrt{\min_i \lambda_i^*}$ lies in U , thus P has radius at least $\frac{1}{d}\ell^*$. On the other hand B has radius ℓ , and we have seen that $P \subseteq B$. Therefore, $\ell \geq \frac{1}{d}\ell^*$. Finally, by our assumption on $\min_i \lambda_i^*$, we have $\min_i \lambda_i^* < 5/\gamma^2$ and so $\ell^* > \gamma/\sqrt{5}$. Therefore, $\ell \geq \gamma/\sqrt{5}d$, which together with (62) guarantees $\frac{\sqrt{2\varepsilon}}{\ell} \leq \frac{\sqrt{5\gamma}}{4} = \frac{1}{2}\sqrt{5\gamma/4}$. Thus, continuing (61):

$$\mathbf{x}^\top P \mathbf{x} \leq (1 - \sqrt{5\gamma/4}) \left(1 + \frac{1}{2}\sqrt{5\gamma/4}\right)^2 \quad (63)$$

Now $(1-x)(1+\frac{x}{2})^2 < 1 - \frac{3}{4}x^2$ for all $x > 0$, thus with $x = \sqrt{5\gamma/4} > \sqrt{\gamma}$ we get:

$$\mathbf{x}^\top P \mathbf{x} < 1 - \frac{3}{4}\gamma \quad (64)$$

By summing (55) and (64), we get:

$$\mathbf{x}^\top P \mathbf{x} + \mathbf{x}^\top Q \mathbf{x} < 1 - \frac{3}{4}\gamma + \frac{\gamma}{4} < 1 \quad (65)$$

Proof of (2). Comparing the eigenvalues of E and E^* , and using $\ell \leq 1$ and $\gamma \leq 1/5$, we obtain:

$$\frac{\lambda_i}{\lambda_i^*} \geq \begin{cases} \frac{(1-\sqrt{5\gamma/4})/\ell^2}{1/\varepsilon^2} \geq \frac{\varepsilon^2}{2} & i \in A_P \\ \varepsilon > \frac{\varepsilon^2}{2} & i \in A_Q \end{cases} \quad (66)$$

Thus the semiaxes lengths of E are at most $\sqrt{2}/\varepsilon$ times those of E^* . Now let E_+^* be the set obtained by scaling E^* by a factor $2\sqrt{2}/\varepsilon = 64\sqrt{2}d^2/\gamma^3$ about its origin $\boldsymbol{\mu}^*$. Note that $\boldsymbol{\mu}^* \in \text{conv}(S_C)$ and, by item (1), $\text{conv}(S_C) \subseteq E$, which implies $\boldsymbol{\mu}^* \in E$. Now, E_+^* contains any set of the form $\mathbf{y} + \frac{1}{2}E_+^*$ if the latter contains $\boldsymbol{\mu}^*$; this includes the set $\frac{\sqrt{2}}{\varepsilon}E^*$ centered in $\boldsymbol{\mu}$, which in turn contains E as we already said.

Proof of (3). We prove that $d(\mathbf{x}, \mathcal{B}_0)^2 < \gamma$ for all $\mathbf{x} \in E$. Since \mathcal{B}_0 is the unit ball, this implies $E \subset \sqrt{1+\gamma}\mathcal{B}_0$. Consider then any such \mathbf{x} . Let again \mathbf{p}, \mathbf{q} be the projections of \mathbf{x} on U and U_\perp respectively. Because $B \subseteq \mathcal{B}_0$, $d(\mathbf{x}, \mathcal{B}_0)^2 \leq d(\mathbf{x}, B)^2 = d(\mathbf{p}, B)^2 + \|\mathbf{q}\|^2$. See again Figure 3, left, but with \mathbf{x} possibly outside \mathcal{B}_0 . For the first term, note that

$$d(\mathbf{p}, B) \leq \max_{i \in A_P} \sqrt{1/\lambda_i} - \ell \quad (67)$$

By definition of λ_i , this yields:

$$d(\mathbf{p}, B)^2 \leq \left(\frac{\ell}{\sqrt{1 - \sqrt{5\gamma/4}}} - \ell \right)^2 \leq \left(\frac{1}{\sqrt{1 - \sqrt{5\gamma/4}}} - 1 \right)^2 \quad (\text{because } \ell \leq 1)$$

Now we show that the right-hand side is bounded by $\frac{3}{4}\gamma$. Consider $f(x) = \frac{1}{\sqrt{1-x}} - 1$ for $x \in [0, 1/2]$. Now $\frac{\partial^2 f}{\partial x^2} = \frac{3}{4}(1-x)^{-5/2} > 0$, so f is convex. Moreover, $f(1/2) = \sqrt{2} - 1 < 0.83 \cdot 1/2$, and clearly $f(0) = 0 \leq 0.83 \cdot 0$. By convexity then, for all $x \in [0, 1/2]$ we have $f(x) \leq 0.83x$ which implies $f(x)^2 < 0.75x^2$. By substituting $x = \sqrt{5\gamma/4}$, for all $\gamma \leq 1/5$ we obtain:

$$d(\mathbf{p}, B)^2 \leq \left(\frac{1}{\sqrt{1 - \sqrt{5\gamma/4}}} - 1 \right)^2 < \frac{3}{4} \cdot \frac{5}{4} \gamma = \frac{15}{16} \gamma \quad (68)$$

Let us now turn to \mathbf{q} . By definition of Q_0 , of Q , and of λ_i for $i \in A_Q$, we have:

$$\|\mathbf{q}\|^2 = \mathbf{x}^\top Q_0 \mathbf{x} \leq \max_{i \in A_Q} \frac{1}{\lambda_i} \mathbf{x}^\top Q \mathbf{x} = \max_{i \in A_Q} \frac{1}{\varepsilon \lambda_i^*} \mathbf{x}^\top Q \mathbf{x} \quad (69)$$

But $\mathbf{x}^\top Q \mathbf{x} \leq 1$ since $\mathbf{x} \in E$, and recalling that $\lambda_i^* \geq 1/\varepsilon^2$ for all $i \in A_Q$, we obtain:

$$\|\mathbf{q}\|^2 \leq \frac{1}{\varepsilon(1/\varepsilon^2)} = \varepsilon < \frac{\gamma}{16} \quad (70)$$

Finally, by summing (68) and (70):

$$d(\mathbf{x}, B_0)^2 \leq d(\mathbf{p}, B)^2 + \|\mathbf{q}\|^2 < \gamma \quad (71)$$

The proof is complete.

3 Supplementary material for Section 6

3.1 Lemma 9

Lemma 9. *Let $b > 0$ be a sufficiently large constant. Let S be a sample of points drawn independently and uniformly at random from X . Let $C = \arg \max_{C_j \in \mathcal{C}} |S \cap C_j|$, let $S_C = S \cap C$, and suppose $|S_C| \geq bd^2 \ln k$. If E is any (possibly degenerate) ellipsoid in \mathbb{R}^d such that $S_C = C \cap E$, then with probability at least $1/2$ we have $|C \cap E| \geq |X| \frac{1}{4k}$. The same holds if we require that $E \cap (S \setminus S_C) = \emptyset$, i.e., that E separates S_C from $S \setminus S_C$.*

Proof. Let $n = |X|$ for short, and for any ellipsoid E let $E_X = E \cap X$. We show that, with C defined as above, (i) with probability at least $1 - 1/4$ we have $|C| \geq n/2k$, and (ii) with probability at least $1 - 1/4$, if $|C| \geq n/2k$ then $|E_X \Delta C| \leq 1/2|C|$ where Δ denotes symmetric difference. By a union bound, then, with probability at least $1/2$ we have $|E \cap C| \geq |C| - |E_X \Delta C| \geq \frac{1}{2}|C| \geq n/4k$.

(i). Let S be the multiset of samples drawn from X , and for every cluster $C_i \in \mathcal{C}$ let N_i be the number of samples in C_i . Let $s = kbd^2 \ln k$; note that $|S| \leq s$ since there are at most k clusters. Now fix any C_i with $|C_i| < \frac{n}{2k}$. Then $\mathbb{E}[N_i] \leq s \frac{|C_i|}{n} < \frac{bd^2 \ln k}{2}$, and by standard concentration bounds (Lemma 4 in this supplementary material), we have $\mathbb{P}(N_i \geq bd^2 \ln k) = \exp(-\Omega(b \ln k))$, which for b large enough drops below $1/4k$. Therefore, the probability that $N_i \geq bd^2 \ln k$ when taking $s \leq kbd^2 \ln k$ samples is at most $1/4k$. By a union bound on all C_i with $|C_i| < n/2k$, then, $|C| \geq n/2k$ with probability $1 - 1/4$.

(ii). Consider now any C_i with $|C_i| \geq n/2k$. We invoke the generalization bounds of Theorem 6 in this supplementary material with $\varepsilon = 1/4k$ and $\delta = 1/4k$, on the hypothesis class \mathcal{H} of all (possibly degenerate) ellipsoids in \mathbb{R}^d . For b large enough, the generalization error of any ellipsoid E that contains S_C is, with probability at least $1 - 1/4k$, at most $1/4k$, which means $|E_X \Delta C_i| \leq n/4k \leq 1/2|C_i|$, as desired. By a union bound on all clusters, with probability at least $1 - 1/4$ this holds for all C_i with $|C_i| \geq n/2k$. The same argument holds if we require E to separate $S \cap C_i$ from $S \setminus C_i$, see again Theorem 6. By a union bound with point (i) above, we have $E \cap C \leq 1/2|C|$ with probability at least $1/2$, as claimed. \square

3.2 Proof of Lemma 3

Let $X_0 = X$ and $N_0 = n$, and for all $i \geq 1$, let X_i be the set of points not yet labeled at the end of round i , let $N_i = |X_i|$, and let $R_i = \mathbb{I}\{N_i \leq N_{i-1}(1 - 1/4k)\}$. Recall that S_C is large enough

so that, by Lemma 9 in this supplementary material, we have $\mathbb{P}(R_i = 1 | X_{i-1}) \geq 1/2$ for all i . For every $t \geq 1$ let $\rho_t = \sum_{i=1}^t R_i$. Note that:

$$N_t \leq N_0(1 - 1/4k)^{\rho_t} < ne^{-\frac{\rho_t}{4k}} \quad (72)$$

If $\rho_t \geq 4k \ln(1/\varepsilon)$, then $N_t < \varepsilon n$ and $\text{RECUR}(X, k, \gamma, \varepsilon)$ stops. The number of rounds executed by $\text{RECUR}(X, k, \gamma, \varepsilon)$ is thus at most $r_\varepsilon = \min\{t : \rho_t \geq 4k \ln(1/\varepsilon)\}$.

Now, for all $i \geq 1$ consider the σ -algebra \mathcal{F}_{i-1} generated by X_0, \dots, X_{i-1} , and define: $Z_i = R_i B_i$, where B_1, B_2, \dots are Bernoulli random variables where each B_i has parameter $1/(2\mathbb{E}[R_i | \mathcal{F}_{i-1}])$. Obviously, $Z_i \leq R_i$, and thus for all t we deterministically have:

$$\rho_t = \sum_{i=1}^t R_i \geq \sum_{i=1}^t Z_i \quad (73)$$

Now note that:

$$\mathbb{E}[Z_i | \mathcal{F}_{i-1}] = \mathbb{E}[R_i | \mathcal{F}_{i-1}] \frac{1}{2\mathbb{E}[R_i | \mathcal{F}_{i-1}]} = \frac{1}{2} \quad (74)$$

Now we can prove the theorem. For the first claim, simply note that $\mathbb{E}[r_\varepsilon] \leq 8k \ln(1/\varepsilon)$, as this is the expected number of fair coin tosses to get $4k \ln(1/\varepsilon)$ heads.

For the second claim, consider any $t \geq 8k \ln n + 6a\sqrt{k} \ln n$. Letting $\zeta_t = \sum_{i=1}^t Z_i$, the event $r_0 \geq t$ implies $\zeta_t < 4k \ln n = \frac{t}{2} - 3a\sqrt{k} \ln n = \mathbb{E}[\zeta_t] - \delta$ where $\delta = 3a\sqrt{k} \ln n$. By Hoeffding's inequality this event has probability at most $e^{-2\delta^2/t}$, and one can check that for all $a \geq 1$ we have $\frac{2\delta^2}{t} \geq a \ln n$.

4 Supplementary material for Section 7

4.1 Proof of Theorem 4

We state and prove two distinct theorems which immediately imply Theorem 4.

Theorem 8. *For all $0 < \gamma < 1/7$, all $d \geq 2$, and every (possibly randomized) learning algorithm, there exists an instance on $n \geq 2\left(\frac{1+\gamma}{8\gamma}\right)^{\frac{d-1}{2}}$ points and $|\mathcal{C}| = 3$ latent clusters such that (1) all clusters have margin γ , and (2) to return with probability $2/3$ a clustering $\hat{\mathcal{C}}$ such that $\Delta(\hat{\mathcal{C}}, \mathcal{C}) = 0$ the algorithm must make $\Omega(n)$ same-cluster queries in expectation.*

Proof. The idea is the following. We define a single set of points $X \subset \mathbb{R}^d$ and randomize over the choice of the latent PSD matrix W ; the claim of the theorem follows by applying Yao's minimax principle. Specifically, we let X be a $\Theta(\sqrt{\gamma})$ -packing of points on the unit sphere in \mathbb{R}^d . We show that, for $\mathbf{x} \in X$ drawn uniformly at random, setting $W = (1 + \gamma) \text{diag}(x_1^2, \dots, x_d^2)$ makes \mathbf{x} an outlier, as its distance $d_W(\mathbf{x}, \mathbf{0})$ from the origin is $1 + \gamma$, while every other point is at distance ≤ 1 . Since there are roughly $(1/\gamma)^d$ such points \mathbf{x} in our set, the bound follows.

We start by defining the points X in terms of their entry-wise squared vectors. Consider $S_d^+ = \mathbb{R}_+^d \cap S_d$ where $S_d = \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\|_2 = 1\}$ is the unit sphere in \mathbb{R}^d . We want to show that there exists a set of $\frac{1}{2}(1/\varepsilon)^{d-1}$ points in S_d^+ whose pairwise distance is bigger than $\varepsilon/2$, where ε will be defined later. To see this, recall that the packing number of the unit ball $B_d = \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\|_2 \leq 1\}$ is $\mathcal{M}(B, \varepsilon) \geq (1/\varepsilon)^d$ —see, e.g., [6]. For $\varepsilon/2$ and $d - 1$, this implies there exists $Y \subseteq B_{d-1}$ such that $|Y| \geq (2/\varepsilon)^{d-1}$ and $\|\mathbf{y} - \mathbf{y}'\|_2 > \varepsilon/2$ for all distinct $\mathbf{y}, \mathbf{y}' \in Y$. Now, consider the lifting function $f : B_{d-1} \rightarrow \mathbb{R}^d$ defined by $f(\mathbf{y}) = (\sqrt{1 - \|\mathbf{y}\|_2^2}, y_1, \dots, y_{d-1})$. Define the lifted set $Z = \{f(\mathbf{y}) : \mathbf{y} \in Y\}$. Clearly, every $\mathbf{z} \in Z$ satisfies $\|\mathbf{z}\|_2 = 1$ and $z_0 \geq 0$, so \mathbf{z} lies on the northern hemisphere of the sphere S_d . Moreover, $\|f(\mathbf{y}) - f(\mathbf{y}')\|_2 \geq \|\mathbf{y} - \mathbf{y}'\|_2$ for any two $\mathbf{y}, \mathbf{y}' \in Y$. Hence, we have a set Z of $(2/\varepsilon)^{d-1}$ points on the d -dimensional sphere such that $\|\mathbf{z} - \mathbf{z}'\|_2 > \varepsilon/2$ for all distinct $\mathbf{z}, \mathbf{z}' \in Z$. But a hemisphere is the union of 2^{d-1} orthants, hence some orthant contains at least $2^{-(d-1)}(2/\varepsilon)^{d-1} = (1/\varepsilon)^{d-1}$ of the points of Z . Without loss of generality we may assume this is the positive orthant and denote the set as Z^+ .

We now define the input set $X \subseteq \mathbb{R}^d$ as follows:

$$X = X^+ \cup X^- = \{\sqrt{z} : z \in Z^+\} \cup \{-\sqrt{z} : z \in Z^+\}$$

Note that $n = |X| = 2|Z^+| = 2(1/\varepsilon)^{d-1}$. Next, we show how every $z \in Z^+$ defines a clustering instance satisfying the constraints of the thesis. For any $z^* \in Z^+$; let $w = (1 + \gamma)z^*$ and $W = \text{diag}(w_1, \dots, w_d)$, which is PSD as required. Define the following three clusters:

$$C' = \{-\sqrt{z^*}\} \quad C'' = \{\sqrt{z^*}\} \quad C = X \setminus (C' \cup C'')$$

where, for $f : \mathbb{R} \rightarrow \mathbb{R}$, $f(\mathbf{x}) = (f(x_1), \dots, f(x_d))$. Since C' and C'' are singletons, they trivially have weak margin γ . We now show that C has weak margin γ w.r.t. to $\mu = \mathbf{0}$; that is, $d_W(\mathbf{x}, \mu)^2 > 1 + \gamma$ for $\mathbf{x} = \pm\sqrt{z^*}$ and $d_W(\mathbf{x}, \mu)^2 \leq 1$ otherwise. First, note that $d_W(\mathbf{x}, \mu)^2 = \langle w, \mathbf{x}^2 \rangle$. Now,

$$d_W(\mathbf{x}, \mu)^2 = \begin{cases} (1 + \gamma) \langle z^*, z^* \rangle = 1 + \gamma & \text{if } \mathbf{x} \in C', C'' \\ (1 + \gamma) \langle z^*, \mathbf{x}^2 \rangle & \text{if } \mathbf{x} \in C \end{cases} \quad (75)$$

However, by construction of Z^+ , we have that for all $\mathbf{x} \in C$ and $z = \mathbf{x}^2$,

$$(\varepsilon/2)^2 \leq \|z - z^*\|_2^2 = \|z\|_2^2 - 2\langle z, z^* \rangle + \|z^*\|_2^2 = 2(1 - \langle z, z^* \rangle)$$

which implies $\langle z^*, \mathbf{x}^2 \rangle \leq 1 - (\varepsilon/2)^2/2 = 1 - \varepsilon^2/8 = 1/(1+\gamma)$ for $\varepsilon = \sqrt{8\gamma/(1+\gamma)}$. Therefore (75) gives $d_W(\mathbf{x}, \mu)^2 = (1 + \gamma) \langle z^*, \mathbf{x}^2 \rangle \leq 1$. This proves C has weak margin γ as desired.

The size of X is:

$$n \geq 2 \left(\frac{1}{\sqrt{8\gamma/(1+\gamma)}} \right)^{d-1} = 2 \left(\frac{1+\gamma}{8\gamma} \right)^{\frac{d-1}{2}}$$

Now the distribution of the instances is defined by taking z^* from the uniform distribution over Z^+ . Consider any deterministic algorithm running over such a distribution. Note that same-cluster queries always return +1 unless at least one of the two queried points is not in C . As C contains all points in X but the symmetric pair $\sqrt{z^*}, -\sqrt{z^*}$ for a randomly drawn z^* , a constant fraction of the points in X must be queried before one element of the pair is found with probability bounded away from zero. Thus, any deterministic algorithm that returns a zero-error clustering with probability at least δ for any constant $\delta > 0$ must perform $\Omega(n)$ queries. By Yao's principle for Monte Carlo algorithms then (see Section 1.4 above), any randomized algorithm that errs with probability at most $\frac{1-\delta}{2} \leq \frac{1}{2}$ for any constant $\delta > 0$ must make $\Omega(n)$ queries as well. \square

Theorem 9. *For all $\gamma > 0$, all $d \geq 48(1 + \gamma)^2$, and every (possibly randomized) learning algorithm, there exists an instance on $n = \Omega(\exp(d/(1 + \gamma)^2))$ points and $|\mathcal{C}| = 2$ latent clusters such that (1) all clusters have margin at least γ , and (2) to return with probability $2/3$ a clustering $\hat{\mathcal{C}}$ such that $\Delta(\hat{\mathcal{C}}, \mathcal{C}) = 0$ the algorithm must make $\Omega(n)$ same-cluster queries in expectation.*

Proof. We exhibit a distribution of instances that gives a lower bound for every algorithm, and then use Yao's minimax principle. Let $p = \frac{1}{2(1+\gamma)}$. Consider a set of vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$ where every entry of each vector $x_{j,i}$ is i.i.d. and it is equal to 1 with probability p . Define $X = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$; note that in general $|X| \leq n$ since the points might not be all distinct. Let $\mathbf{x}^* = \mathbf{x}_n$, $C = \{\mathbf{x}_1, \dots, \mathbf{x}_{n-1}\}$, $C' = \{\mathbf{x}^*\}$. The latent clustering is $\mathcal{C} = \{C, C'\}$, and the matrix and center of \mathcal{C} are respectively $W = \text{diag}(\mathbf{x}^*)$ and $\mathbf{c} = \mathbf{0}$. The algorithms receive in input a random permutation of X ; clearly, if it makes $o(|X|)$ queries, then it has vanishing probability to find \mathbf{x}^* , which is necessary to return the latent clustering \mathcal{C} .

Now we claim that, if $d \geq 48(1 + \gamma)^2$, then we can set $n = \Omega\left(\exp\left(\frac{d}{48(1+\gamma)^2}\right)\right)$ and with constant probability we will have (i) $|X| = \Omega(n)$, and (ii) C, C' have margin γ . This is sufficient, since the theorem then follows by applying Yao's minimax principle.

Let us first bound the probability that $|X| < n$. Note that for any two points $\mathbf{x}_i, \mathbf{x}_{i'}$ with $i \neq i'$ we have $\mathbb{P}(\mathbf{x}_i = \mathbf{x}_{i'}) = ((1-p)^2 + p^2)^d < (1 - \frac{1}{2(1+\gamma)})^d < e^{-\frac{d}{2(1+\gamma)}}$. Therefore, by a simple union bound over all pairs, $\mathbb{P}(|X| < n) < n^2 e^{-\frac{d}{2(1+\gamma)}}$.

Next, we want show that, loosely speaking, $d_W(\mathbf{x}, \mathbf{c})^2 \simeq dp$ for $\mathbf{x} \in C'$ whereas $d_W(\mathbf{x}, \mathbf{c})^2 \simeq dp^2$ for $\mathbf{x} \in C$; this will give the margin.

Now, for any \mathbf{x} ,

$$d_W(\mathbf{x}, \mathbf{c})^2 = \sum_{i=1}^d x_i^* (x_i - 0)^2 = \begin{cases} \sum_{i=1}^d x_i^* x_i \sim B(d, p^2) & \mathbf{x} \in C \\ \sum_{i=1}^d x_i^* \sim B(d, p) & \mathbf{x} \in C' \end{cases} \quad (76)$$

Where in the last equality we use the fact that the entries are unary, and where with the notation $B(d, p)$ we refer to a vector of length d where each entry is equal to 1 with probability p . Let $\mu = dp^2$ and $\mu' = dp$, let $\varepsilon = 1/(1+\sqrt{2})$, and define

$$\phi = \mu(1 + \varepsilon), \quad \phi' = \mu'(1 - \varepsilon\sqrt{p}) \quad (77)$$

By standard tail bounds,

$$\mathbb{P}(d_W(\mathbf{x}, \mathbf{c})^2 \geq \phi) \leq e^{-\frac{\varepsilon^2 \mu}{3}} \quad \text{for } \mathbf{x} \in C \quad (78)$$

$$\mathbb{P}(d_W(\mathbf{x}, \mathbf{c})^2 < \phi') < e^{-\frac{\varepsilon^2 p \mu'}{3}} = e^{-\frac{\varepsilon^2 \mu}{3}} \quad \text{for } \mathbf{x} \in C' \quad (79)$$

By a union bound on all points, the margin γ_C of C fails to satisfy the following inequality with probability at most $|X|e^{-\frac{\varepsilon^2 \mu}{3}} \leq ne^{-\frac{\varepsilon^2 \mu}{3}}$:

$$1 + \gamma_C = \frac{\min_{\mathbf{x} \notin C} d_W(\mathbf{x}, \mathbf{c})^2}{\max_{\mathbf{x} \in C} d_W(\mathbf{x}, \mathbf{c})^2} \geq \frac{\phi'}{\phi} = \frac{dp(1 - \varepsilon\sqrt{p})}{dp^2(1 + \varepsilon)} = \frac{1 - \varepsilon\sqrt{p}}{p(1 + \varepsilon)} \geq \frac{1}{2p} = 1 + \gamma \quad (80)$$

where the penultimate inequality holds since $\frac{1 - \varepsilon\sqrt{p}}{1 + \varepsilon} \geq \frac{1}{2}$ for our values of p and ε . Note that, since $p = \frac{1}{2(1+\gamma)}$ and $n \leq \frac{1}{c} \exp\left(\frac{d}{48(1+\gamma)^2}\right) + 1$,

$$ne^{-\frac{\varepsilon^2 \mu}{3}} = ne^{-\frac{dp^2}{12}} = ne^{-\frac{d}{48(1+\gamma)^2}} \quad (81)$$

By one last union bound, the probability that $|X| = n$ and $\gamma_C \geq \gamma$ is at least

$$1 - ne^{-\frac{d}{48(1+\gamma)^2}} - n^2 e^{-\frac{d}{2(1+\gamma)}} \quad (82)$$

If $d \geq \frac{48}{(1+\gamma)^2}$, then we can let $n = \Omega\left(e^{\frac{d}{48(1+\gamma)^2}}\right)$ while ensuring the above probability is bounded away from 0.

The rest of the proof and the application of Yao's principle is essentially identical to the proof of Theorem 8 above. \square

5 Comparison with SCQ- k -means

In this section we compare our algorithm to SCQ- k -means of [1]. We show that, in our setting, SCQ- k -means fails even on very simple instances, although it can still work under (restrictive) assumptions on γ , W , and the centers.

SCQ- k -means works as follows. First, the center of mass μ_C of some cluster C is estimated using $\mathcal{O}(\text{poly}(k, 1/\gamma))$ SCQ queries; second, all points in X are sorted by their distance from μ_C and the radius of C is found via binary search. The binary search is done using same-cluster queries between the sorted points and any point already known to be in C . The margin condition ensures that, if we have an accurate enough estimate of μ_C , then the binary search will be successful (there are no inversions of the sorted points w.r.t. their cluster). This approach thus yields a $\mathcal{O}(\ln n)$ SCQ queries bound (the number of queries to estimate μ_C is independent of n).

It is easy to see that this algorithm relies crucially on (1) each cluster C must be spherical, and (2) the center of the sphere must coincide with the centroid μ_C . In formal terms, the setting of [1] is a special cases of ours where for all C we have $W_C = I_d$ and $\mathbf{c} = \mathbb{E}_{\mathbf{x} \in C}[\mathbf{x}]$. If any of these two assumptions does not hold, then it is easy to construct instances where [1] fails to recover the clusters and, in fact, achieves error very close to a completely random labeling. Formally:

Lemma 10. For any fixed $d \geq 2$, any $p \in (0, 1)$, and any sufficiently small $\gamma > 0$, there are arbitrarily large instances on n points and $k = 2$ clusters on which SCQ- k -means incurs error $\Delta(\hat{\mathcal{C}}, \mathcal{C}) \geq \frac{1-p}{2}$ with probability at least $1 - p$.

Sketch of the proof. We describe the generic instance on n points for $d = 2$. The latent clustering \mathcal{C} is formed by two clusters C_1, C_2 of size respectively $n_1 = n \frac{1+p}{2}$ and $n_2 = n \frac{1-p}{2}$. In C_1 , half of the points are in $(1, 0)$ and half in $(-1, 0)$. In C_2 , all points are in $(0, \frac{\sqrt{1+\gamma}}{2})$. (One can in fact perturb the instance so that all points are distinct without impairing the proof). For both clusters, the center coincide with their center of mass, $\mu_1 = (0, 0)$ and $\mu_2 = (0, \frac{\sqrt{1+\gamma}}{2})$. For both clusters, the latent metric is given by the PSD matrix $W = \begin{pmatrix} .25 & 0 \\ 0 & 1 \end{pmatrix}$. It is easy to see that $d_W(\mathbf{x}, \mu_1)^2 = 1/4$ if $\mathbf{x} \in C_1$ and $d_W(\mathbf{x}, \mu_1)^2 = (1+\gamma)/4$ if $\mathbf{x} \in C_2$, and so C_1 has margin exactly γ . On the other hand C_2 has margin γ since $d_W(\mathbf{x}, \mu_2)^2 = 0$ if $\mathbf{x} \in C_2$ and $d_W(\mathbf{x}, \mu_2)^2 > 0$ otherwise.

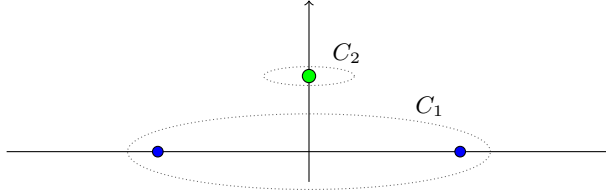


Figure 4: A bad instance for SCQ- k -means. With good probability, the algorithm classifies all points in a single cluster, incurring error $\simeq 1/2$, the same as a random labeling.

Now consider SCQ- k -means. The algorithm starts by sampling at least $\frac{k \ln(k)}{\gamma^4}$ points from X and setting $\hat{\mu}$ to the average of the points with the majority label. By standard concentration bounds then, for γ small enough, with probability at least $1 - p$ the majority cluster will be C_1 and the estimate $\hat{\mu}$ of its center of mass $(0, 0)$ will be sufficiently close to μ_1 that the ordering of all points in X by their Euclidean distance w.r.t. $\hat{\mu}$ will set all of C_2 before all of C_1 . But since $n_2 = n \frac{1-p}{2}$, the median of the sorted sequence will be a point of C_1 . Thus the binary search will make its first query on a point of C_1 and will continue thereafter classifying all of X as belonging to C_1 . Thus the algorithm will output the clustering $\hat{\mathcal{C}} = \{X, \emptyset\}$ which gives $\Delta(\hat{\mathcal{C}}, \mathcal{C}) = \frac{1-p}{2}$. \square

Next, we show that the approach [1] still works if one relaxes the assumption $W = I$, at the price of strengthening the margin γ . Let λ_{\max} and $\lambda_{\min} > 0$ be, respectively, the largest and smallest eigenvalues of W . The condition number κ_W of W is the ratio $\lambda_{\max}/\lambda_{\min}$. If κ_W is not too large, then W does not significantly alter the Euclidean metric, and the ordering of the points is preserved. Formally:

Lemma 11. Let κ_W be the condition number of W . If every cluster C has margin at least $\kappa_W - 1$ with respect to its center of mass μ_C , and if we know μ_C , then we can recover C with $\mathcal{O}(\ln n)$ SCQ queries.

Proof. Fix any cluster C and let $\mu = \mu_C$. For any $\mathbf{z} \in \mathbb{R}^d$ we have $\lambda_{\min} \|\mathbf{z}\|_2^2 \leq \|\mathbf{z}\|_W^2 \leq \lambda_{\max} \|\mathbf{z}\|_2^2$ where λ_{\min} and λ_{\max} are, respectively, the smallest and largest eigenvalue of W . Sort all other points \mathbf{x} by their Euclidean distance $\|\mathbf{x} - \mu\|_2$ from μ . Then, for any $\mathbf{x} \in C$ and any $\mathbf{y} \notin C$ we have:

$$\frac{\|\mathbf{y} - \mu\|_2^2}{\|\mathbf{x} - \mu\|_2^2} \geq \frac{\lambda_{\min} \|\mathbf{y} - \mu\|_W^2}{\lambda_{\max} \|\mathbf{x} - \mu\|_W^2} = \frac{1}{\kappa_W} \frac{d(\mathbf{y}, \mu)^2}{d(\mathbf{x}, \mu)^2} > \frac{1 + \gamma}{\kappa_W} \quad (83)$$

Hence, if $\gamma \geq \kappa_W - 1$, there is $r \geq 0$ such that $\|\mathbf{x} - \mu\|_2 \leq r$ for all $\mathbf{x} \in C$ and $\|\mathbf{y} - \mu\|_2 \geq r$ all $\mathbf{y} \notin C$. We can thus recover C via binary search as in [1]. \square

As a final remark, we observe that the above approach is rather brittle, since κ_W is unknown (because W is), and if the condition $\kappa_W \leq 1 + \gamma$ fails, then once again the binary search can return a clustering far from the correct one.

6 Comparison with metric learning

In this section we show that metric learning, a common approach to latent cluster recovery and related problems, does not solve our problem even when combined with same-cluster and comparison queries. Intuitively, we want to learn an approximate distance \hat{d} that preserves the ordering of the distances between the points. That is, for all $\mathbf{x}, \mathbf{y}, \mathbf{z} \in X$, $d(\mathbf{x}, \mathbf{y}) \leq d(\mathbf{x}, \mathbf{z})$ implies $\hat{d}(\mathbf{x}, \mathbf{y}) \leq \hat{d}(\mathbf{x}, \mathbf{z})$. If this holds then d and \hat{d} are equivalent from the point of view of binary search. To simplify the task, we may equip the algorithm with an additional *comparison query* CMP, which takes in input two pairs of points \mathbf{x}, \mathbf{x}' and \mathbf{y}, \mathbf{y}' from X and tells precisely whether $d(\mathbf{x}, \mathbf{x}') \leq d(\mathbf{y}, \mathbf{y}')$ or not. It turns out that, even with SCQ+CMP queries, learning such a \hat{d} requires to query essentially all the input points.

Theorem 10. *For any $d \geq 3$, learning any \hat{d} such that, for all $\mathbf{x}, \mathbf{y}, \mathbf{z} \in X$, if $d(\mathbf{x}, \mathbf{y}) \leq d(\mathbf{x}, \mathbf{z})$ then $\hat{d}(\mathbf{x}, \mathbf{y}) \leq \hat{d}(\mathbf{x}, \mathbf{z})$, requires $\Omega(n)$ SCQ+CMP queries in the worst case, even with an arbitrarily large margin γ .*

Proof. We reduce the problem of learning the order of pairwise distances induced by W , which we call ORD, to the problem of learning a separator hyperplane, which we call SEP and whose query complexity is linear in n .

Problem SEP is as follows. The inputs are a set $X = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subset \mathbb{R}^d$ (the observations) and a set $\mathcal{H} = \{\mathbf{h}_1, \dots, \mathbf{h}_k\} \subset \mathbb{R}_+^d$ (the hypotheses). We require that $\mathbf{h}_j \in \mathbb{R}_+^d$. We have oracle access to $\sigma : X \rightarrow \{+1, -1\}$ such that $\sigma(\cdot) = \text{sgn}\langle \mathbf{h}, \cdot \rangle$ for some $\mathbf{h} \in \mathcal{H}$. The output is the $\mathbf{h} \in \mathcal{H}$ that agrees with σ . We assume \mathcal{H}, X support a margin: $\exists \varepsilon > 0$, possibly dependent on the instance, such that $\text{sgn}\langle \mathbf{h}, \mathbf{x} \rangle = \text{sgn}\langle \mathbf{h}, \mathbf{x}' \rangle$ for all \mathbf{x}' with $\|\mathbf{x} - \mathbf{x}'\| \leq \varepsilon$. (Note that this is *not* the cluster margin γ).

Let $Q_{\text{ORD}}(n)$ and $Q_{\text{SEP}}(n)$ be the query complexities of ORD and SEP on n points. We show:

Lemma 12. $Q_{\text{ORD}}(3n) \leq Q_{\text{SEP}}(n)$.

Proof. Let $X = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subseteq \mathbb{R}^d$ be the input points for SEP and let $\mathbf{h} \in \mathbb{R}_+^d$ be the target hypothesis. By scaling the dataset we can assume $\|\mathbf{x}_i\| \leq \varepsilon$ for any desired ε (even dependent on n). We define an instance of ORD on $n' = 3n$ points as follows. First, $W = \text{diag}(\mathbf{h})$. Second, the input set is $X' = S_1 \cup \dots \cup S_n$ where for $i = 1, \dots, n$ we define $S_i = \{\mathbf{a}_i, \mathbf{b}_i, \mathbf{c}_i\}$ with:

$$\mathbf{a}_i = 6^i \cdot \mathbf{1} \tag{84}$$

$$\mathbf{b}_i = 2 \cdot \mathbf{a}_i \tag{85}$$

$$\mathbf{c}_i = 3 \cdot \mathbf{a}_i + \mathbf{x}_i \tag{86}$$

We first show that a solution to ORD gives a solution of SEP. Suppose indeed that for all pairs of points $\{\mathbf{q}, \mathbf{p}\}, \{\mathbf{x}, \mathbf{y}\}$ we know whether $d_W(\mathbf{q}, \mathbf{p}) \leq d_W(\mathbf{x}, \mathbf{y})$. This is equivalent to knowing the output of $\text{CMP}(\{\mathbf{q}, \mathbf{p}\}, \{\mathbf{x}, \mathbf{y}\})$, which is

$$\text{CMP}(\{\mathbf{q}, \mathbf{p}\}, \{\mathbf{x}, \mathbf{y}\}) = \text{sgn} \langle \mathbf{h}, (\mathbf{q} - \mathbf{p})^2 - (\mathbf{x} - \mathbf{y})^2 \rangle \tag{87}$$

Consider then the point $\mathbf{q} = \mathbf{c}_i, \mathbf{p} = \mathbf{x} = \mathbf{b}_i, \mathbf{y} = \mathbf{a}_i$ for each i . Then:

$$\text{CMP}(\{\mathbf{q}, \mathbf{p}\}, \{\mathbf{x}, \mathbf{y}\}) = \text{sgn} \langle \mathbf{h}, (\mathbf{a}_i - \mathbf{b}_i)^2 - (\mathbf{b}_i - \mathbf{c}_i)^2 \rangle \tag{88}$$

$$= \text{sgn} \langle \mathbf{h}, (\mathbf{a}_i)^2 - (-\mathbf{a}_i - \mathbf{x}_i)^2 \rangle \tag{89}$$

$$= \text{sgn} \langle \mathbf{h}, 2 \cdot 6^i \mathbf{x}_i - \mathbf{x}_i^2 \rangle \tag{90}$$

$$= \text{sgn} \left\langle \mathbf{h}, \mathbf{x}_i \left(1 - \frac{\mathbf{x}_i}{2 \cdot 6^i} \right) \right\rangle \tag{91}$$

By the margin hypothesis, for ε small enough this equals $\text{sgn}(\langle \mathbf{h}, \mathbf{x}_i \rangle)$, i.e., the label of \mathbf{x}_i in SEP.

We now show that all the other queries reveal no information about the solution of SEP. Suppose then the points are not in the form $\mathbf{q} = \mathbf{c}_i, \mathbf{p} = \mathbf{x} = \mathbf{b}_i, \mathbf{y} = \mathbf{a}_i$. Without loss of generality, we can assume that $\mathbf{q} > \mathbf{p}$ and $\mathbf{q} \geq \mathbf{x} > \mathbf{y}$. It is then easy to see that, for ε small enough, $(\mathbf{q} - \mathbf{p})^2 - (\mathbf{x} - \mathbf{y})^2 > 0$ or $(\mathbf{q} - \mathbf{p})^2 - (\mathbf{x} - \mathbf{y})^2 < 0$. This holds independently of the \mathbf{x}_i and of W and therefore gives no information about the solution of SEP.

It follows that, if we can solve ORD in $f(3n)$ CMP queries, then we can solve SEP in $f(n)$ queries. Finally, note that adding SCQ queries does not reduce the query complexity (e.g., let X lie in a single cluster). For the same reason, we can even assume an arbitrarily large cluster margin γ . \square

It remains to show that SEP requires $\Omega(n)$ CMP queries in the worst case. This is well known, but we need to ensure that $\mathcal{H} \subset \mathbb{R}_+^d$ and that any $h \in \mathcal{H}$ supports a margin as described above.

Consider the following set $X = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subseteq \mathbb{R}^3$:

$$\mathbf{x}_i = (1 - \delta, -\cos(\theta_i), -\sin(\theta_i)) \quad (92)$$

where $\theta_i = i \frac{\pi}{2n}$ and δ is sufficiently small. Let $\mathcal{H} = \{\mathbf{h}_1, \dots, \mathbf{h}_n\}$, where

$$\mathbf{h}_j = (1, \cos(\theta_j), \sin(\theta_j)) \quad (93)$$

Note that $\mathcal{H} \subset \mathbb{R}_+^d$ as required. Clearly:

$$\langle \mathbf{h}_j, \mathbf{x}_i \rangle = \begin{cases} -\delta & \text{if } j = i \\ 1 - (\delta + \cos(\theta_i - \theta_j)) & \text{if } j \neq i \end{cases} \quad (94)$$

By choosing $\delta = \frac{1 - \cos(\pi/2n)}{2}$ we have $\text{sgn} \langle \mathbf{h}, \mathbf{x}_i \rangle = -1$ if and only if $i = j$. Clearly, any algorithm needs to probe $\Omega(n)$ labels to learn h with constant probability for some $h \in \mathcal{H}$. Finally, note that any h supports a margin, as required. \square

References

- [1] Hassan Ashtiani, Shrinu Kushagra, and Shai Ben-David. Clustering with same-cluster queries. In *Advances in Neural Information Processing Systems 29*, pages 3216–3224. 2016.
- [2] Devdatt Dubhashi and Alessandro Panconesi. *Concentration of Measure for the Analysis of Randomized Algorithms*. Cambridge University Press, New York, NY, USA, 1st edition, 2009.
- [3] Hervé Fournier and Olivier Teytaud. Lower bounds for comparison based evolution strategies using vc-dimension and sign patterns. *Algorithmica*, 59(3):387–408, March 2011.
- [4] Rajeev Motwani and Prabhakar Raghavan. *Randomized Algorithms*. Cambridge University Press, USA, 1995.
- [5] Shai Shalev-Shwartz and Shai Ben-David. *Understanding Machine Learning: From Theory to Algorithms*. Cambridge University Press, USA, 2014.
- [6] Roman Vershynin. *High-Dimensional Probability: An Introduction with Applications in Data Science*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 2018.