# Preclustering Algorithms for Imprecise Points

Document status and date:
Published: 01/06/2022

Document Version:
Publisher's PDF, also known as Version of Record (includes final page, issue and volume numbers)

Please check the document version of this publication:

• A submitted manuscript is the version of the article upon submission and before peer-review. There can be
important differences between the submitted version and the official published version of record. People
interested in the research are advised to contact the author for the final version of the publication, or visit the
DOI to the publisher's website.
• The final author version and the galley proof are versions of the publication after peer review.
• The final published version features the final layout of the paper including the volume, issue and page
numbers.

Link to publication

# Preclustering Algorithms for Imprecise Points

**Mohammad Ali Abam[1]** (iD) · **Mark de Berg[2]** (iD) · **Sina Farahzad[1]** · **Mir Omid Haji Mirsadeghi[3]** · **Morteza Saghafian[3]**

## Abstract

We study the problem of *preclustering* a set $B$ of imprecise points in $\mathbb{R}^d$: we wish to cluster the regions specifying the potential locations of the points such that, no matter where the points are located within their regions, the resulting clustering approximates the optimal clustering for those locations. We consider $k$-center, $k$-median, and $k$-means clustering, and obtain the following results. Let $B := \{b_1, \ldots, b_n\}$ be a collection of disjoint balls in $\mathbb{R}^d$, where each ball $b_i$ specifies the possible locations of an input point $p_i$. A partition $\mathcal{C}$ of $B$ into subsets is called an $(f(k), \alpha)$-preclustering (with respect to the specific $k$-clustering variant under consideration) if (i) $\mathcal{C}$ consists of $f(k)$ preclusters, and (ii) for any realization $P$ of the points $p_i$ inside their respective balls, the cost of the clustering on $P$ induced by $\mathcal{C}$ is at most $\alpha$ times the cost of an optimal $k$-clustering on $P$. We call $f(k)$ the *size* of the preclustering and we call $\alpha$ its *approximation ratio*. We prove that, even in $\mathbb{R}^1$, one may need at least $3k - 3$ preclusters to obtain a bounded approximation ratio—this holds for the $k$-center, the $k$-median, and the $k$-means problem—and we present a $(3k, 1)$ preclustering for the

✉ Mohammad Ali Abam
  abam@sharif.edu

  Mark de Berg
  m.t.d.berg@tue.nl

  Sina Farahzad
  farahzad@ce.sharif.edu

  Mir Omid Haji Mirsadeghi
  mirsadeghi@sharif.edu

  Morteza Saghafian
  morteza.saghafian65@student.sharif.edu

[1]  Department of Computer Engineering, Sharif University of Technology, Tehran, Iran

[2]  Department of Mathematics and Computer Science, TU Eindhoven, Eindhoven, The Netherlands

[3]  Department of Mathematical Sciences, Sharif University of Technology, Tehran, Iran

$k$-center problem in $\mathbb{R}^1$. We also present various preclusterings for balls in $\mathbb{R}^d$ with $d \geq 2$, including a $(3k, \alpha)$-preclustering with $\alpha \approx 13.9$ for the $k$-center and the $k$-median problem, and $\alpha \approx 193.9$ for the $k$-means problem.

**Keywords** Computational Geometry · Clustering · Imprecise Points

## 1 Introduction

Clustering is one of the most important and widely studied problems in unsupervised learning. It comes in many different flavors, depending on the type of data to be clustered, the measure used to assess the quality of a clustering, and so on. In this paper we are interested in geometric clustering, where the data are points in $\mathbb{R}^d$, and we consider three well-known centroid-based clustering methods, namely $k$-center, $k$-median, and $k$-means, on so-called imprecise points.

In (the geometric version of) centroid-based clustering one is given a set $P$ of $n$ points in $\mathbb{R}^d$, where $d$ is a fixed constant, and an integer $k$. The goal is to partition $P$ into $k$ subsets $P_1, \ldots, P_k$ and assign a centroid $q_i$ to each cluster $P_i$ such that the cost of the resulting clustering is minimized. In the $k$-center problem the cost of the clustering is defined as $\max_{1 \leq i \leq k} \max_{p \in P_i} |pq_i|$, where $|pq|$ denotes the Euclidean distance between two points $p$ and $q$. In the $k$-median problem the cost of a clustering is defined as $\sum_{1 \leq i \leq k} \sum_{p \in P_i} |pq_i|$, and in the $k$-means problem it is defined as $\sum_{1 \leq i \leq k} \sum_{p \in P_i} |pq_i|^2$. Given a collection of centroids it is always optimal to define the clusters by assigning each point in $P$ to its nearest centroid. Thus an equivalent definition of the $k$-center problem, for instance, is to find a collection of $\{q_1, \ldots, q_k\}$ as centroids that minimizes $\max_{p \in P} \min_{1 \leq i \leq k} |pq_i|$. In other words, we want to find $k$ congruent balls of minimum radius that together cover all points in $P$.

The $k$-center problem in $\mathbb{R}^d$ is NP-hard for $d \geq 2$ when $k$ is part of the input [17]. For the Euclidean $k$-center problem a PTAS exists, as shown by Agarwal and Procopiuc [1]. (For the $k$-center problem in general metric spaces, a PTAS does not exists; for this case an $r$-approximation algorithm with $r < 2$ is not possible unless P=NP, and several 2-approximation algorithms are known [8,21].) The $k$-median and $k$-means problems are also NP-hard for $d \geq 2$ [16,17], and they admit a PTAS as well [2,6,9,12].

In the traditional setting the locations of the input points are known exactly. In practice this may not always be the case: typically locations are measured using GPS or other devices that are not completely accurate, or the points may move around inside a given region. This leads to the study of geometric algorithms on so-called *imprecise points*. Here, instead of specifying the exact coordinates of each input point, we specify a region for each point where it may be located. For points in the plane the regions are typically disks or squares. Over the past decade, many problems have been studied for imprecise points, including convex hulls (compute the smallest (or largest) possible convex hull of a set of imprecise points [10,15]), Delaunay triangulations (preprocess a set of imprecise points such that for any given instantiation of the points in the given regions we can compute the Delaunay triangulation quickly [4]), separability problems [20], and more [13,14,18].

As already mentioned, the regions where each point can be located are typically considered to be disks. A common assumption is that these disks are pairwise disjoint. In other words, one assumes that the imprecision in the locations of the points is relatively small compared to the inter-point distances. Note that if one would allow all disks to have a non-empty common intersection, then they do not give any information about the relative positions of the points. Hence, for most problems one cannot obtain interesting results when the disks are allowed to intersect. This is also the case for our problem (as we remark below, after defining the problem more precisely) and so we will assume that the disks (or balls, in $\mathbb{R}^d$) that specify the locations are pairwise disjoint.

One may try to overcome the disjointness condition by using a statistical model, for example where the potential location of each point is given by a Gaussian distribution. See [7,11] for some examples of papers studying the clustering problem in such a setting. A probabilistic setting by itself does not solve the problem, however: if all distributions are essentially the same—for example, all Gaussians are centered at (almost) the same point—then the distributions do not give any information about the relative positions of the points. Hence, additional assumptions are also needed here. Also note that the goal of a probabilistic analysis (including smoothed analysis [22]) is, in some sense, to get rid of worst-case examples. We, on the other hand, want to deal with worst-case placement of the points inside their imprecision region.

We remark that, besides the notion of imprecise points that we use, there has also been work on so-called *uncertain points*. Here each point has a finite set of possible locations (rather than a region where the point may lie), each with an associated probability. A special case is the existential model, where there is only one location whose probability is smaller than 1, so that with a certain probability the point is not present at all. The uncertain-points model typically leads to substantially different questions and results from the imprecise-points model, because the former has probabilities associated with the locations (so it is closer to the statistical models mentioned above) and because of its discrete nature.

Finally, we mention the notion of perturbation resilience introduced by Bilu and Linial [3]. In the context of clustering, this notions means that one can change the inter-point distances by at most a given factor, without changing the optimal solution; see for example [5] and the references therein. This is somewhat related to our model, since the disks specifying the possible point locations can be seen as regions in which an adversary is allowed to move a point. However, we do *not* assume that these disks are such that the optimal solution does not change—for this disjointness is far from sufficient—and thus consider different questions from the ones studied in the papers dealing with perturbation resilience.

*Problem statement and notation* In this paper we study the $k$-center, $k$-median, and $k$-means problem for imprecise points. The input is a set $B := \{b_1, \ldots, b_n\}$ of disjoint (closed) balls in $\mathbb{R}^d$, each representing the possible locations of an input point. Our goal is to compute a *preclustering* of the imprecise points, that is, a partition of $B$ into a collection $\mathcal{C}$ of subsets called *preclusters* that gives a good clustering for any possible realization of the points inside the input balls. Next we define this more formally.

For a (precise) point set $P$, let $\text{OPT}_\infty(P, k)$ denote the cost of an optimal $k$-center clustering on $P$, that is,

$$\text{OPT}_\infty(P, k) := \min_{q_1,\ldots,q_k \in \mathbb{R}^d} \max_{p \in P} \min_{1 \le i \le k} |pq_i|.$$

The cost of an optimal solution for the $k$-median and $k$-means problem on a set $P$ are denoted by $\text{OPT}_1(P, k)$ and $\text{OPT}_2(P, k)$, respectively.[1] Now consider an imprecise point set specified by a set $B = \{b_1, \ldots, b_n\}$ of balls. A point set $P := \{p_1, \ldots, p_n\}$ such that $p_i \in b_i$ for all $1 \le i \le n$ is called a $B$-*instance*. A preclustering $\mathcal{C}$ of the set $B$ into preclusters $B_i$ induces a clustering on any $B$-instance $P$ in a natural manner, namely by creating a cluster $P_i := \{p \in P : p \in B_i\}$ for every precluster $B_i \in \mathcal{C}$. The cost of the preclustering $\mathcal{C}$ on $P$, denoted by $\mathcal{C}\text{-COST}_\infty(P)$ for the $k$-center problem, is defined as the cost of the induced clustering on $P$ if we choose the centroid of each cluster $P_i$ optimally, namely by solving the 1-clustering problem on $P_i$. (We use the term *preclustering* since we already perform the partitioning into clusters before we know the actual (precise) location of the points. The computation of the optimal centers for each of the clusters $P_i$, which we can only do after the actual locations are known, can be considered the postprocessing step of the clustering.) So for the $k$-center problem we have

$$\mathcal{C}\text{-COST}_\infty(P) := \max_{B_i \in \mathcal{C}} \min_{q \in \mathbb{R}^d} \max_{p \in P_i} |pq|.$$

The preclustering costs for the $k$-median and $k$-means problem are denoted by $\mathcal{C}\text{-COST}_1(P)$ and $\mathcal{C}\text{-COST}_2(P)$, respectively, and they are defined similarly. To quantify the quality of a preclustering $\mathcal{C}$ on $B$ (with respect to the $k$-clustering problem under consideration) we define $\mathcal{C}$ to be a $(f(k), \alpha)$-*preclustering* if

– $\mathcal{C}$ consists of $f(k)$ preclusters,
– $\mathcal{C}\text{-COST}(P) \le \alpha \cdot \text{OPT}(P, k)$ for any $B$-instance $P$.

We call $f(k)$ the *size* of the preclustering and we call $\alpha$ its *approximation ratio*. Observe that if the balls in $B$ would have a non-empty common intersection, then any preclustering with fewer than $n$ preclusters may have an arbitrarily bad approximation ratio, even for the 2-center problem. This is the reason that we assume (as mentioned earlier and as is often done in papers on imprecise points) that the balls in $B$ are disjoint.

*Our results* Ideally, we would like to have a $(k, 1)$-preclustering, but this is typically impossible. This leads to the question: what is the smallest value for $f(k)$ such that we can always obtain an $(f(k), 1)$-preclustering? More generally, which trade-offs are possible between the size $f(k)$ of the preclustering and its approximation ratio $\alpha$?

In Sect. 2 we study this problem in $\mathbb{R}^1$. We show that there are input sets $B$ that require at least $3k - 3$ preclusters to get a bounded approximation ratio; this holds

---

[1] The subscript $\infty$ in $\text{OPT}_\infty$ refers to the fact that if $d_i$ denotes the distance of point $p_i \in P$ to its nearest center, then we are minimizing the norm of the vector $\langle d_1, \ldots, d_n \rangle$ in the $\ell_\infty$-metric. For $k$-median and $k$-means we are minimizing the norm in the $\ell_1$-metric and in the squared $\ell_2$-metric, respectively.

**Fig. 1** Illustration of the lower-bound construction for $k = 5$: a collection of $k - 1$ groups of three intervals (in grey), each group consisting of a left and right interval of length 1 separated by a gap of length $\varepsilon$, and a middle interval inside this gap. The points in the $B$-instance used in the proof are shown slightly above the intervals for clarity

for the $k$-center problem, the $k$-median problem, as well as the $k$-means problem. We complement this result by proving that any set $B$ of intervals in $\mathbb{R}^1$ admits a $(3k, 1)$-preclustering for the $k$-center problem. This preclustering can be computed in polynomial time.

In Sect. 3 we consider the $d$-dimensional version of the problem for $d \geq 2$. We give an example showing that here a $(3k, 1)$-preclustering does not always exist, and we present a $(3k, \alpha)$-preclustering with $\alpha \approx 13.9$ for the $k$-center and $k$-median problem, and $\alpha \approx 193.9$ for the $k$-means problem. A different parameterization of the strategy gives a $(6k, 3)$-preclustering for $k$-center and $k$-median, and a $(6k, 9)$-preclustering for $k$-means in $\mathbb{R}^2$.

Because we allow $f(k) > k$ it may also be possible to obtain an approximation factor $\alpha < 1$. In Sect. 4, we study the question whether we can achieve any approximation factor $\varepsilon > 0$ by choosing $f(k)$ large enough or not. We first prove the following negative result. It is not always possible to obtain any given approximation ratio $\varepsilon > 0$ for the $k$-median and the $k$-means problem, with a preclustering size $f(k)$ that is independent from $n$. After that we obtain tight asymptotic bounds on the size of the preclustering needed to obtain any given approximation ratio $\varepsilon > 0$ for the $k$-center problem. In particular, we prove that $\Theta(\lceil \sqrt{d}/\varepsilon \rceil^d \cdot k)$ preclusters are always sufficient and sometimes $\Omega(\lceil 1/(\varepsilon\sqrt{d}) \rceil^d \cdot k)$ preclusters are necessary to obtain an approximation ratio $\varepsilon$.

## 2 The 1-Dimensional Problem

We begin by proving that even in $\mathbb{R}^1$—here the input balls are disjoint intervals on the line—preclusterings with only $k$ preclusters cannot always guarantee a good approximation ratio. In fact, we sometimes need as much as $3k - 3$ preclusters in any preclustering with bounded approximation ratio.

**Theorem 1** *For any integer $k \geq 2$ and any given $\alpha$, there is a set $B$ of disjoint intervals in $\mathbb{R}^1$ that does not admit a $(k', \alpha)$-preclustering with $k' < 3k - 3$. This holds for $k$-center, $k$-median, as well as $k$-means clustering.*

**Proof** Let $B$ be a collection of $3k - 3$ disjoint intervals in $\mathbb{R}^1$ consisting of $k - 1$ groups of three intervals each. The left and right interval in each group have length 1 and are at distance $\varepsilon$ from each other, where $\varepsilon$ is a sufficiently small number that will be specified later. The middle interval from the group lies in the gap between the left and right interval with its center at the center of the gap; see Fig. 1.

Now consider a preclustering $\mathcal{C} = \{B_1, \ldots, B_{k'}\}$. If $k' < 3k - 3$, then there is at least one precluster containing two consecutive intervals, $b_i$ and $b_j$. Assume without

loss of generality that length($b_i$) $\geq$ length($b_j$), and consider the $B$-instance in which each point $p_t$ is placed in its interval $b_t \in B$ as follows.

- If $t = i$ or $b_t$ is a middle interval, then $p_t$ lies at the center of $b_t$.
- If $t \neq i$ and $b_t$ is a left interval, then $p_t$ lies at the right endpoint of $b_t$.
- If $t \neq i$ and $b_t$ is a right interval, then $p_t$ lies at the left endpoint of $b_t$.

Note that with this placement we have $|p_i p_j| \geq 1/2$. We will argue that by choosing $\varepsilon$ appropriately we get the desired result.

First consider the $k$-center problem. Note that $\text{OPT}_\infty(P, k) \leq \varepsilon/2$. Indeed, by putting a centroid at the center of each of the $k-1$ gaps and one centroid at $p_i$, all points in $P$ are at distance at most $\varepsilon/2$ from a centroid. On the other hand, $\mathcal{C}\text{-COST}_\infty(P) \geq 1/4$ since the centroid for the cluster containing $p_i$ and $p_j$ is at distance at least $1/4$ from $p_i$ or $p_j$. Hence,

$$\frac{\mathcal{C}\text{-COST}_\infty(P)}{\text{OPT}_\infty(P, k)} \geq \frac{1/4}{\varepsilon/2} = \frac{1}{2\varepsilon}.$$

For $\varepsilon < 1/(2\alpha)$ we thus enforce an approximation ratio greater than $\alpha$.

The argument for $k$-median and $k$-means is similar. For $k$-median we have $\text{OPT}_1(P, k) \leq 2(k-1)(\varepsilon/2)$ and $\mathcal{C}\text{-COST}_1(P) \geq 1/2$, so $\varepsilon < 1/(2(k-1)\alpha)$ enforces an approximation ratio greater than $\alpha$, while for $k$-means we have $\text{OPT}_2(P, k) \leq 2(k-1)(\varepsilon/2)^2$ and $\mathcal{C}\text{-COST}_2(P) \geq 2(1/4)^2$, so it suffices to have $\varepsilon < \sqrt{1/(4(k-1)\alpha)}$.                                                                $\square$

**Remark 1** The construction in the proof of Theorem 1 uses an input set $B$ of size $3k-3$. We can easily generate an input set with the same behavior for any $n \geq 3k - 3$, by adding another $n - 3k + 3$ tiny intervals inside one of the gaps between a left and a right interval from the same group.

Theorem 1 states that for some problem instances any preclustering with fewer than $3k-3$ preclusters has arbitrarily large approximation ratio. We now show how to obtain a 1-approximation with only $3k$ preclusters for the $k$-center problem. We assume from now on that $n > 3k$, otherwise we can trivially create a zero-cost solution with at most $3k$ preclusters.

Before we describe our preclustering strategy, we first generalize the $k$-center problem in $\mathbb{R}^1$ from points to intervals. In this generalization the input is a collection $B$ of $n$ disjoint intervals, and the goal is to find a collection $\mathcal{I} := \{I_1, \ldots, I_k\}$ of intervals that together cover all intervals in $B$ and such that the maximum radius of the intervals in $\mathcal{I}$ is minimized. (The radius of an interval is half its length.) We denote the value of an optimal solution $\mathcal{I}$ to the $k$-center problem on $B$ by $\text{OPT}_\infty(B, k)$, so $\text{OPT}_\infty(B, k) := \max_{I_i \in \mathcal{I}} \text{radius}(I_i)$.

Our preclustering algorithm is now as follows.

PRECLUSTERING-1D$(B, k)$

1. Sort the intervals in $B$ by radius, such that $\text{radius}(b_1) \geq \cdots \geq \text{radius}(b_n)$.
2. For each $k' \in \{0, \ldots, 2k\}$ do the following.

(a) Let $\{B_1, \ldots, B_{(3k-k')}\}$ be an optimal $(3k-k')$-center clustering on $\{b_{k'+1}, \ldots, b_n\}$, and let $\text{OPT}_\infty(\{b_{k'+1}, \ldots, b_n\}, 3k - k')$ be its cost.

(b) Let $\mathcal{C}(k')$ be the preclustering $\{\{b_1\}, \ldots, \{b_{k'}\}, B_1, \ldots, B_{(3k-k')}\}$.

3. Of all preclusterings $\mathcal{C}(0), \ldots, \mathcal{C}(2k)$ found in Step 2, let $\mathcal{C}(k')$ be the one that minimizes $\text{OPT}_\infty(\{b_{k'+1}, \ldots, b_n\}, 3k - k')$. Let $\mathcal{C} := \mathcal{C}(k')$ and return $\mathcal{C}$.

**Theorem 2** *Let $B$ be a set of disjoint intervals in $\mathbb{R}^1$. Then $B$ admits a $(3k, 1)$-preclustering for the $k$-center problem and such a preclustering can be computed in polynomial time.*

**Proof** Obviously PRECLUSTERING-1D$(B, k)$ gives a preclustering $\mathcal{C}$ with $3k$ preclusters. Next we prove that $\mathcal{C}$ has approximation ratio 1. Let $P$ be a $B$-instance, and let $Q \in \{q_1, \ldots, q_k\}$ be an optimal set of centroids for the $k$-center problem on $P$. Thus by placing an interval of radius $\text{OPT}_\infty(P, k)$ centered at each centroid $q_i \in Q$, we cover all points in $P$. By assigning each point in $P$ to its nearest centroid in $Q$, with ties broken arbitrarily, we obtain a partition of $P$ into $k$ clusters. This partition induces a preclustering $\mathcal{C}^*$ of size $k$ on $B$. We use $\mathcal{C}^*$ to define two types of intervals: *outer intervals*, which are the leftmost or rightmost interval in any of the preclusters $B_i \in \mathcal{C}^*$, and *inner intervals*, which are the remaining intervals. Note that the number of outer intervals is at most $2k$. Define $k^*$ as the largest $k'$ such that $b_1, \ldots, b_{k'}$ are all outer intervals, where $b_1, \ldots, b_n$ is the sorted set of intervals obtained in Step 1 of the algorithm. Since $b_{k^*+1}$ is an inner interval, we have

$$\text{OPT}_\infty(P, k) \geq \text{radius}(b_{k^*+1}). \tag{1}$$

The preclustering $\mathcal{C} := \mathcal{C}(k')$ returned by our algorithm minimizes the value of $\text{OPT}_\infty(\{b_{k'+1}, \ldots, b_n\}, 3k - k')$. Note that

$$\mathcal{C}\text{-COST}_\infty(P) \leq \text{OPT}_\infty(\{b_{k'+1}, \ldots, b_n\}, 3k - k'),$$

since the intervals $b_1, \ldots, b_{k'}$ are all in singleton preclusters and an interval covering all intervals in a precluster $B_i$ obviously covers all points from $P$ in those interval. Hence,

$$\mathcal{C}\text{-COST}_\infty(P) \leq \text{OPT}_\infty(\{b_{k^*+1}, \ldots, b_n\}, 3k - k^*).$$

It remains to argue that $\text{OPT}_\infty(P, k) \geq \text{OPT}_\infty(\{b_{k^*+1}, \ldots, b_n\}, 3k - k^*)$. To this end, we create a collection $\mathcal{I}$ of intervals as follows.

– For each outer interval $b_j$ with $j > k^*$ we create an interval equal to $b_j$.
– For each precluster $B_i \in \mathcal{C}^*$ that has at least one inner interval, we create a minimum-length interval covering all inner intervals of $B_i$.

Note that $\mathcal{I}$ contains at most $3k - k^*$ intervals, and that these intervals together cover all intervals in $\{b_{k^*+1}, \ldots, b_n\}$. Hence,

$$\max_{I \in \mathcal{I}} \text{radius}(I) \geq \text{OPT}_\infty(\{b_{k^*+1}, \ldots, b_n\}, 3k - k^*).$$
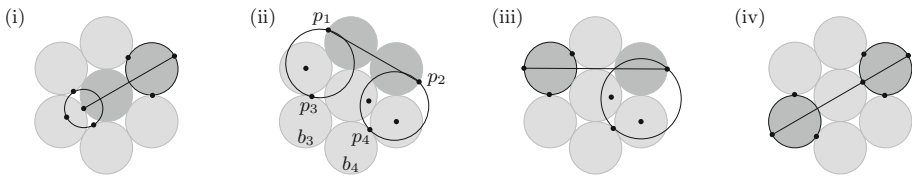
**Fig. 2** The seven balls shown in the figure do not admit a $(3k, 1)$-preclustering for $k = 2$

Moreover, $\text{OPT}_\infty(P, k) \geq \text{radius}(I)$ for any $I \in \mathcal{I}$. Indeed, if $I$ is equal to an outer interval $b_j$ with $j > k^*$ then $\text{OPT}_\infty(P, k) \geq \text{radius}(b_j)$ by Inequality (1), and otherwise $I$ is the minimum-length interval covering all inner intervals of some precluster $B_i$. (In the latter case we also have $\text{OPT}_\infty(P, k) \geq \text{radius}(I)$ because in any $B$-instance the cluster of $B_I$ includes a point in both outer intervals) We conclude that

$$\text{OPT}_\infty(P, k) \geq \max_{I \in \mathcal{I}} \text{radius}(I) \geq \text{OPT}_\infty(\{b_{k^*+1}, \ldots, b_n\}, 3k - k^*).$$

It remains to argue that $\text{PRECLUSTERING-1D}(B, k)$ can be implemented to run in polynomial time. The most time-consuming step is Step 2a, which can be implemented to run in $O(n^2 k)$ time using dynamic programming in a straightforward manner. □

Theorem 2 only holds for the $k$-center problem. In the next section we present a more general algorithm, which not only works in higher dimensions but also for $k$-median and $k$-means. The size of the computed preclustering will not be as good as what is provided by Theorem 2.

## 3 The *d*-Dimensional Problem

In the previous section we saw that for some problem instances any preclustering with fewer than $3k - 3$ preclusters has an arbitrarily large approximation ratio. The result is stated for $\mathbb{R}^1$ but it also holds in $\mathbb{R}^d$ for $d > 1$: we can use exactly the same construction, replacing the intervals by $d$-dimensional balls whose centers lie on the $x_1$-axis. We also presented an algorithm giving a $(3k, 1)$-preclustering for intervals in $\mathbb{R}^1$, for the $k$-center problem.

Figure 2 shows that a $(3k, 1)$-preclustering is not always possible for the $k$-center problem in $\mathbb{R}^2$. The figure shows a set $B$ of seven unit balls, with one central ball touching the other six balls. For $k = 2$ a preclustering of size $3k$ would use five singleton preclusters and one precluster with two balls. There are four combinatorially distinct ways of choosing the precluster of two balls, indicated by the dark grey balls in parts (i)–(iv) of the figure. For each case, a $B$-instance is shown (the black dots), and the optimal solution to the 2-center problem for the instance is shown (the two black circles). The best preclustering is the one in part (ii). Here the two points $p_1, p_2$ in the dark grey balls are placed at distance 4 from each other, so $\mathcal{C}\text{-COST}_\infty(P) = 2$. The point $p_3$ inside the ball $b_3$ is placed as close to $p_1$ as possible, while $p_4$ is placed as close to $p_2$ as possible. The other points are placed such that they are either

contained in the ball with diameter $p_1 p_3$ or in the ball with diameter $p_2 p_4$. Hence, $\text{OPT}_\infty(P) = (\sqrt{13} - 1)/2$. The balls in this construction are not disjoint, but we can slightly scale them to obtain an instance where any $(3k, \alpha)$-preclustering has $\alpha \geq 2/((\sqrt{13} - 1)/2) - \varepsilon \approx 1.54$.

We now present a preclustering strategy that works for $k$-center, $k$-means and $k$-median in any dimension. It is similar to, and actually somewhat simpler than, the preclustering algorithm we presented for the 1-dimensional $k$-center problem.

PRECLUSTERING-DD$(B, k)$

1. Sort the balls in $B$ by radius, such that $\text{radius}(b_1) \geq \cdots \geq \text{radius}(b_n)$.
2. Define $B_{\text{small}} := \{b_{2k+1}, \ldots, b_n\}$; we call the balls in $B_{\text{small}}$ *small*. Let $\{P_1, \ldots, P_k\}$ be an optimal $k$-center (or $k$-median, or $k$-means) clustering on the point set $\text{centers}(B_{\text{small}}) := \{c_j : 2k + 1 \leq j \leq n\}$, where $c_j$ is the center of the ball $b_j$. Let $\{B_1, \ldots, B_k\}$ be the preclustering on $B_{\text{small}}$ induced by it.
3. Return the preclustering $\mathcal{C} := \{\{b_1\}, \ldots, \{b_{2k}\}, B_1, \ldots, B_k\}$.

Before we analyze the algorithm's approximation ratio, we note that, depending on the dimension $d$ and the value of $k$, we may not be able to implement Step 2 efficiently [16, 17]. However, instead of computing an optimal $k$-clustering on the centers of the small balls, we can also compute a $(1 + \varepsilon')$-approximation of the optimal clustering. For an appropriate $\varepsilon' = O(\varepsilon)$ this increases the approximation ratio by only a factor $1 + \varepsilon$, as explained later.

Obviously PRECLUSTERING-DD$(B, k)$ gives a preclustering of size $3k$. To analyze the approximation ratio, we use the following lemma.

**Lemma 3** *For any B-instance P the output* $\mathcal{C} := \{\{b_1\}, \ldots, \{b_{2k}\}, B_1, \ldots, B_k\}$ *of the algorithm satisfies:*

*(i)* $\mathcal{C}\text{-COST}_\infty(P) \leq \text{OPT}_\infty(P, k) + 2 \cdot \text{radius}(b_{2k+1})$
*(ii)* $\mathcal{C}\text{-COST}_1(P) \leq \text{OPT}_1(P, k) + 2 \sum_{j=2k+1}^{n} \text{radius}(b_j)$
*(iii)* $\sqrt{\mathcal{C}\text{-COST}_2(P)} \leq \sqrt{\text{OPT}_2(P, k)} + 2\sqrt{\sum_{j=2k+1}^{n} \text{radius}(b_j)^2}.$

**Proof** We first prove part (i) of the lemma. Let $P$ be any $B$-instance, let $p_j \in P$ denote the point inside $b_j$, and let $c_j$ be the center of $b_j$. Recall that $P_i \subset P$ is the subset of points in the instance corresponding to the precluster $B_i$. Define $P_{\text{small}} := \{p_{2k+1}, \ldots, p_n\}$ to be the set of points from $P$ in the small balls, and define $C_{\text{small}} := \{c_{2k+1}, \ldots, c_n\}$. Note that $P_{\text{small}} = P_1 \cup \cdots \cup P_k$ and that

$$|p_j c_j| \leq \text{radius}(b_j) \leq \text{radius}(b_{2k+1}) \tag{2}$$

for all $p_j \in P_{\text{small}}$. We define the following sets of centroids:

– Let $Q := \{q_1, \ldots, q_k\}$ be the set of centroids in an optimal $k$-center solution for the entire point set $P$. We have

$$\max_{p_j \in P_{\text{small}}} \min_{q_i \in Q} |p_j q_i| \leq \max_{p_j \in P} \min_{q_i \in Q} |p_j q_i| = \text{OPT}_\infty(P, k). \tag{3}$$

- Let $Q' := \{q'_1, \ldots, q'_k\}$ be the set of centroids in the optimal $k$-center clustering on $C_{\text{small}}$ used in Step 2 of the algorithm. Thus

$$\max_{c_i \in C_{\text{small}}} \min_{q'_j \in Q'} |c_i q'_j| = \text{OPT}_\infty(C_{\text{small}}, k) \leq \max_{c_i \in C_{\text{small}}} \min_{q_j \in Q} |c_i q_j|. \tag{4}$$

- Let $Q'' := \{q''_1, \ldots, q''_k\}$, where $q''_i$ is the optimal centroid for $P_i$. Note that for all $P_i$ we have

$$\max_{p_j \in P_i} |p_j q''_i| \leq \max_{p_j \in P_i} |p_j q'_i|. \tag{5}$$

Since the total cost of the singleton preclusters is trivially zero, we have

$\mathcal{C}\text{-}\text{COST}_\infty(P)$
$= \max_{1 \leq i \leq k} \max_{p_j \in P_i} |p_j q''_i|$
$\leq \max_{1 \leq i \leq k} \max_{p_j \in P_i} |p_j q'_i|$ (Inequality (5))
$\leq \max_{1 \leq i \leq k} \max_{p_j \in P_i} (|p_j c_j| + |c_j q'_i|)$ (triangle inequality)
$\leq \text{radius}(b_{2k+1}) + \max_{1 \leq i \leq k} \max_{p_j \in P_i} |c_j q'_i|$ (Inequality (2))
$\leq \text{radius}(b_{2k+1}) + \max_{c_j \in C_{\text{small}}} \min_{q'_i \in Q'} |c_j q'_i|$ (definition of $Q'$)
$\leq \text{radius}(b_{2k+1}) + \max_{c_j \in C_{\text{small}}} \min_{q_i \in Q} |c_j q_i|$ (Inequality (4))
$\leq \text{radius}(b_{2k+1}) + \max_{p_j \in P_{\text{small}}} \min_{q_i \in Q} (|c_j p_j| + |p_j q_i|)$ (triangle inequality)
$\leq 2 \cdot \text{radius}(b_{2k+1}) + \max_{p_j \in P_{\text{small}}} \min_{q_i \in Q} |p_j q_i|$ (Inequality (2))
$\leq 2 \cdot \text{radius}(b_{2k+1}) + \text{OPT}_\infty(P, k)$ (Inequality (3))

To prove part (ii) of the lemma, which deals with the $k$-median problem, we define the sets $Q$, $Q'$ and $Q''$ as above, but now for the $k$-median problem. Note that Inequality (2) still holds while Inequalities (3)–(5) hold if we replace the max-operator by a summation. Part (ii) can thus be derived using a similar derivation as for part (i).

To prove part (iii), which deals with the $k$-means problem, we need to work with squared distances. We define $Q, Q', Q''$ as above, but now for the $k$-means problem. Note that Inequality (2) still holds, while Inequalities (3)–(5) hold if we replace the max-operator with a summation and all distance values with their squared values. We denote these variants by Inequalities (3*)−(5*). For squared distances the triangle inequality does not hold. Instead we use the triangle inequality for $\ell_2$-metric, which is called Minkowsky Inequality. A similar computation as above can now be used to prove part (iii); we have

$\sqrt{\mathcal{C}\text{-}\text{COST}_2(P)}$
$= \sqrt{\sum_{i=1}^k \sum_{p_j \in P_i} |p_j q''_i|^2}$
$\leq \sqrt{\sum_{i=1}^k \sum_{p_j \in P_i} |p_j q'_i|^2}$ (Inequality (5))
$\leq \sqrt{\sum_{i=1}^k \sum_{p_j \in P_i} (|p_j c_j| + |c_j q'_i|)^2}$ (triangle inequality)
$\leq \sqrt{\sum_{i=1}^k \sum_{p_j \in P_i} |p_j c_j|^2} + \sqrt{\sum_{i=1}^k \sum_{p_j \in P_i} |c_j q'_i|^2}$ (Minkowsky inequality)
$\leq \sqrt{\sum_{j=2k+1}^n \text{radius}(b_j)^2} + \sqrt{\sum_{i=1}^k \sum_{p_j \in P_i} |c_j q'_i|^2}$ (Inequality (2))

On the other hand

$$\sqrt{\sum_{i=1}^{k}\sum_{p_j\in P_i}|c_jq_i'|^2}$$
$$\leq \sqrt{\sum_{c_j\in C_{\text{small}}}\min_{q_i'\in Q'}|c_jq_i'|^2} \qquad \text{(definition of } Q')$$
$$\leq \sqrt{\sum_{c_j\in C_{\text{small}}}\min_{q_i\in Q}|c_jq_i|^2} \qquad \text{(Inequality (4))}$$
$$\leq \sqrt{\sum_{c_j\in C_{\text{small}}}\min_{q_i\in Q}\left(|c_jp_j|+|p_jq_i|\right)^2} \qquad \text{(triangle inequality)}$$
$$\leq \sqrt{\sum_{c_j\in C_{\text{small}}}|c_jp_j|^2}+\sqrt{\sum_{c_j\in C_{\text{small}}}\min_{q_i\in Q}|p_jq_i|^2} \qquad \text{(Minkowsky inequality)}$$
$$\leq \sqrt{\sum_{j=2k+1}^{n}\text{radius}(b_j)^2}+\sqrt{\sum_{c_j\in C_{\text{small}}}\min_{q_i\in Q}|p_jq_i|^2} \qquad \text{(Inequality (2))}$$
$$\leq \sqrt{\sum_{j=2k+1}^{n}\text{radius}(b_j)^2}+\sqrt{\text{OPT}_2(P,k)} \qquad \text{(Inequality (3))}$$

By adding up the above inequalities we conclude part (iii) of the lemma. $\square$

**Remark 2** Note that algorithm PRECLUSTERING-DD computes an optimal clustering on centers($B_{\text{small}}$), which may be expensive. Instead we can also work with a $(1+\varepsilon)$-approximation to an optimal clustering, for some $\varepsilon > 0$. Then the statement of Lemma 3 becomes

(i) $\mathcal{C}\text{-COST}_\infty(P) \leq (1+\varepsilon)\cdot\text{OPT}_\infty(P,k)+(2+\varepsilon)\cdot\text{radius}(b_{2k+1})$
(ii) $\mathcal{C}\text{-COST}_1(P) \leq (1+\varepsilon)\cdot\text{OPT}_1(P,k)+(2+\varepsilon)\cdot\sum_{j=2k+1}^{n}\text{radius}(b_j)$
(iii) $\sqrt{\mathcal{C}\text{-COST}_2(P)} \leq \sqrt{1+\varepsilon}\cdot\sqrt{\text{OPT}_2(P,k)}+(1+\sqrt{1+\varepsilon})\cdot\sqrt{\sum_{j=2k+1}^{n}\text{radius}(b_j)^2}$.

The proof of this modified version is the same as before. The only changes will appear in Inequality (4), where we get an extra multiplicative factor $1+\varepsilon$ for $k$-center and $k$-median, and $\sqrt{1+\varepsilon}$ for $k$-means. We will use this modified version of Lemma 3 later, to prove the second part of Theorem 6.

The lemma above shows that our preclustering gives an additive error that depends on the radii of the small balls. The following two lemmas will be used to turn this into a multiplicative error.

For the next two lemmas, we define $r_d^*$ as the smallest possible radius of any ball that can intersect three disjoint unit balls in $\mathbb{R}^d$ (for instance, see Fig. 3). More formally, we define

$$r_d^* := \inf\{\text{radius}(b) : b \text{ is a ball that intersects three disjoint unit balls in } \mathbb{R}^d\}.$$

**Lemma 4** *We have*

*(i)* $\text{OPT}_\infty(P,k) \geq r_d^*\cdot\text{radius}(b_{2k+1})$
*(ii)* $\text{OPT}_1(P,k) \geq r_d^*\cdot\sum_{j=2k+1}^{n}\text{radius}(b_j)$
*(iii)* $\text{OPT}_2(P,k) \geq (r_d^*)^2\cdot\sum_{j=2k+1}^{n}\text{radius}(b_j)^2$

**Proof** For part (i) notice that by the Pigeonhole Principle an optimal clustering must have a cluster containing at least three points from $\{p_1, \ldots, p_{2k+1}\}$. The cost of this cluster is lower bounded by the radius of the smallest ball intersecting three balls of radius at least $b_{2k+1}$, which is in turn lower bounded by $r_d^* \cdot \mathrm{radius}(b_{2k+1})$.

For part (ii) let $P_1, P_2, \ldots, P_k$ be the clusters in an optimal $k$-median clustering on $P$, and let $q_i$ be the centroid of $P_i$ in this clustering. Let $B_i$ be the set of balls corresponding to the points in $P_i$. We claim that

$$\sum_{p_j \in P_i} |p_j q_i| \geq r_d^* \cdot \left( \left( \sum_{b_j \in B_i} \mathrm{radius}(b_j) \right) - s_i \right). \tag{6}$$

where $s_i$ is the sum of the radii of the two largest balls in $B_i$. To show this, let $b(q_i, r)$ be the ball of radius $r$ centered at $q_i$, and let $P_i(r) := \{p_j \in P_i : b_j \cap b(q_i, r) \neq \emptyset\}$ be the set of points in $P_i$ whose associated ball intersects $b(q_i, r)$, and let $B_i(r)$ be the set of their respective balls. Since for sufficiently large $r$ we have $P_i = P_i(r)$, it suffices to show that for all $r > 0$ we have

$$\sum_{p_j \in P_i(r)} |p_j q_i| \geq r_d^* \cdot \left( \left( \sum_{p_j \in B_i(r)} \mathrm{radius}(b_j) \right) - s_i(r) \right).$$

where $s_i(r)$ is the sum of the radii of the two largest balls in $B_i(r)$. To prove this, consider this inequality as $r$ increases from $r = 0$ to $r = \infty$. As long as $|P_i(0)| \leq 2$ the right-hand side is zero and so the inequality is obviously true. As we increase $r$ further, $b(q_i, r)$ starts intersecting more and more balls from $B_i$. Consider what happens to the inequality when $b(q_i, r)$ starts intersecting another ball $b_\ell \in B_i$. Then $p_\ell$ is added to $P_i(r)$, so the left-hand side of the inequality increases by $|p_\ell q_i|$, which is at least $r$. The right-hand side increases by at most $r_d^*$ times the radius of the third-largest ball in $B_i(r)$. By definition of $r_d^*$, if three balls intersect a ball of radius $r$ then the smallest has radius at most $r/r_d^*$. Hence, the right-hand side increases by at most $r$ and the inequality remains true.
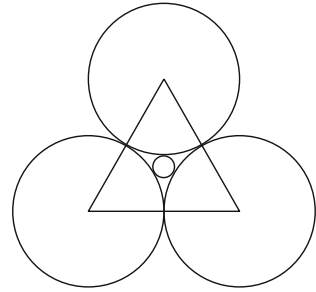
Recall that $b_1, \ldots, b_{2k}$ are the $2k$ largest balls in $B$. Hence, summing Inequality (6) over all clusters $P_1, \ldots, P_k$ gives

$$\begin{aligned}
\mathrm{OPT}_1(P, k) &= \sum_{i=1}^{k} \sum_{p_j \in P_i} |p_j q_i| \\
&\geq r_d^* \cdot \left( \sum_{i=1}^{k} \sum_{b_j \in B_i} \mathrm{radius}(b_j) - \sum_{j=1}^{2k} \mathrm{radius}(b_j) \right) \\
&= r_d^* \cdot \sum_{j=2k+1}^{n} \mathrm{radius}(b_j).
\end{aligned}$$

For part (iii) we define $P_1, P_2, \cdots, P_k$ as above, but now for the $k$-means problem. Also we define $q_i$ as the centroid of $P_i$ and $B_i$ as the set of balls corresponding to the points in $P_i$. Now the same proof as (ii) works if we replace all distances with squared distances. In fact, the inequality (6) becomes

$$\sum_{p_j \in P_i} |p_j q_i|^2 \geq (r_d^*)^2 \cdot \left( \left( \sum_{b_j \in B_i} \mathrm{radius}(b_j)^2 \right) - s_i \right).$$

**Fig. 3** The figure shows the smallest possible ball intersecting three disjoint unit balls in 2D. The larger balls are the unit balls and the radius of the small ball is $r_2^* = \frac{2}{\sqrt{3}} - 1$

where $s_i$ is the sum of the squared radii of the two largest balls in $B_i$.

Therefore, we obtain

$$
\begin{aligned}
\mathrm{OPT}_2(P, k) &= \sum_{i=1}^{k} \sum_{p_j \in P_i} |p_j q_i|^2 \\
&\geq (r_d^*)^2 \cdot \left( \sum_{i=1}^{k} \sum_{b_j \in B_i} \mathrm{radius}(b_j)^2 - \sum_{j=1}^{2k} \mathrm{radius}(b_j)^2 \right) \\
&= (r_d^*)^2 \cdot \sum_{j=2k+1}^{n} \mathrm{radius}(b_j)^2.
\end{aligned}
$$

□

**Lemma 5** *For all $d \geq 2$ we have $r_d^* = \frac{2}{\sqrt{3}} - 1$.*

**Proof** It is easy to see that $r_d^* \leq r_2^*$, since any configuration of three disjoint unit disks in the plane, with a fourth disk intersecting all three, can be extended to $\mathbb{R}^d$ by embedding the centers of the balls on a 2-dimensional plane in $\mathbb{R}^d$. Next we show that $r_d^* \geq r_2^*$ for all $d \geq 2$, which implies that $r_d^* = r_2^*$.

Let $d \geq 2$ and let $b, b', b''$ be three disjoint unit balls in $\mathbb{R}^d$. Let $c, c', c''$ denote the centers of $b, b',$ and $b''$, respectively, and let $h$ be a 2-dimensional plane containing $c, c', c''$. Let $D$ be a smallest ball that intersects $b, b', b''$ and whose center is restricted to lie on $h$. Then $\mathrm{radius}(D) \geq r_2^*$. We claim that $D$ is in fact a smallest ball intersecting $b, b', b''$ even if we do not restrict the center of this ball to be on $h$. Indeed, if a ball $D'$ with center $q \notin h$ intersects $b, b', b''$, then the ball of the same radius as $D'$ and whose center is the orthogonal projection of $q$ onto $h$ also intersects $b, b', b''$.

It remains to show that $r_2^* = \frac{2}{\sqrt{3}} - 1$. The configuration minimizing the radius of the smallest ball intersecting $b, b', b''$ is where $b, b', b''$ are pairwise touching, resulting in the claimed bound—see Fig. 3. □

We are now ready to prove the following theorem.

**Theorem 6** *Let $B$ be a set of disjoint balls in $\mathbb{R}^d$ with $d \geq 2$. Then*

*(i) there exists a $(3k, 7+4\sqrt{3})$-preclustering for the $k$-center and the $k$-median problem,*

*(ii) there exists a $(3k, 97+56\sqrt{3})$-preclustering for the $k$-means problem.*

*Moreover, $(3k, 7+4\sqrt{3}+\varepsilon)$-preclusterings for the $k$-center problem and for the $k$-median problem, can be computed in polynomial time. For the $k$-means problem, we can compute a $(3k, 97+56\sqrt{3}+\varepsilon)$-preclustering in polynomial time, assuming $d$ is a constant.*

**Proof** Parts (i) and (ii) follow immediately by putting together Lemmas 3–5. More precisely, for the $k$-center problem we have

$$
\begin{aligned}
\mathcal{C}\text{-}\mathrm{COST}_\infty(P) &\leq \mathrm{OPT}_\infty(P, k) + 2 \cdot \mathrm{radius}(b_{2k+1}) \\
&\leq \left(1 + \tfrac{2}{r_d^*}\right) \cdot \mathrm{OPT}_\infty(P, k) = (7 + 4\sqrt{3}) \cdot \mathrm{OPT}_\infty(P, k).
\end{aligned}
$$

For the $k$-median problem we have

$$
\begin{aligned}
\mathcal{C}\text{-}\mathrm{COST}_1(P) &\leq \mathrm{OPT}_1(P, k) + 2\textstyle\sum_{j=2k+1}^n \mathrm{radius}(b_j) \\
&\leq \left(1 + \tfrac{2}{r_d^*}\right) \cdot \mathrm{OPT}_1(P, k) = (7 + 4\sqrt{3}) \cdot \mathrm{OPT}_1(P, k).
\end{aligned}
$$

For the $k$-means problem we have

$$
\begin{aligned}
\sqrt{\mathcal{C}\text{-}\mathrm{COST}_2(P)} &\leq \sqrt{\mathrm{OPT}_2(P, k)} + 2\sqrt{\textstyle\sum_{j=2k+1}^n \mathrm{radius}(b_j)^2} \\
&\leq \left(1 + \tfrac{2}{r_d^*}\right) \cdot \sqrt{\mathrm{OPT}_2(P, k)}
\end{aligned}
$$

and therefore,

$$
\mathcal{C}\text{-}\mathrm{COST}_2(P) \leq \left(1 + \frac{2}{r_d^*}\right)^2 \cdot \mathrm{OPT}_2(P, k) = (97 + 56\sqrt{3}) \cdot \mathrm{OPT}_2(P, k).
$$

It remains to argue that we can compute a preclustering whose approximation ratio is as claimed in polynomial time. Recall that both $k$-center and $k$-median, and $k$-means for constant $d$, admit a PTAS [1,2,6,9,12,19], that is, for any given $\varepsilon' > 0$ we can compute a $(1 + \varepsilon')$-approximation to an optimal clustering in polynomial time. To obtain the result, in Step 2 of PRECLUSTERING-DD$(B, k)$ we compute a $(1 + \varepsilon')$-approximation of the optimal clustering for appropriate $\varepsilon'$ that will be introduced later. The resulting algorithm runs in polynomial time. Then the statement of Lemma 3 becomes (see Remark 2 above):

(i) $\mathcal{C}\text{-}\mathrm{COST}_\infty(P) \leq (1 + \varepsilon') \cdot \mathrm{OPT}_\infty(P, k) + (2 + \varepsilon') \cdot \mathrm{radius}(b_{2k+1})$
(ii) $\mathcal{C}\text{-}\mathrm{COST}_1(P) \leq (1 + \varepsilon') \cdot \mathrm{OPT}_1(P, k) + (2 + \varepsilon') \cdot \sum_{j=2k+1}^n \mathrm{radius}(b_j)$
(iii) $\sqrt{\mathcal{C}\text{-}\mathrm{COST}_2(P)} \leq \sqrt{1 + \varepsilon'} \cdot \sqrt{\mathrm{OPT}_2(P, k)} + (1 + \sqrt{1 + \varepsilon'}) \cdot \sqrt{\sum_{j=2k+1}^n \mathrm{radius}(b_j)^2}.$

Using inequalities of Lemma 4 we get

(i) $\mathcal{C}\text{-}\mathrm{COST}_\infty(P) \leq (1 + \varepsilon' + \tfrac{2+\varepsilon'}{r_d^*}) \cdot \mathrm{OPT}_\infty(P, k)$
(ii) $\mathcal{C}\text{-}\mathrm{COST}_1(P) \leq (1 + \varepsilon' + \tfrac{2+\varepsilon'}{r_d^*}) \cdot \mathrm{OPT}_1(P, k)$
(iii) $\mathcal{C}\text{-}\mathrm{COST}_2(P) \leq \left((\sqrt{1 + \varepsilon'} + \tfrac{1+\sqrt{1+\varepsilon'}}{r_d^*})\right)^2 \cdot \mathrm{OPT}_2(P, k).$

By taking $\varepsilon' := \varepsilon / (1 + \tfrac{1}{r_d^*})$ for the $k$-center and $k$-median problem the approximation ratio for the whole algorithm will increase by $\varepsilon$.

For the $k$-means problem, given a positive $\varepsilon$, we need an $\varepsilon'$ such that

$$\left(\left(\sqrt{1+\varepsilon'}+\frac{1+\sqrt{1+\varepsilon'}}{r_d^*}\right)\right)^2 \leq \left(1+\frac{2}{r_d^*}\right)^2 + \varepsilon$$

Note that the left-hand side goes to $(1+\frac{2}{r_d^*})^2$ as $\varepsilon' \to 0$, so by choosing $\varepsilon'$ sufficiently small, as a function of $\varepsilon$ we can satisfy the inequality.                                  □

*Generalizing the solution* We can generalize the above theorem in order to control the number of preclusters for various approximations. Let $r_d^p$ be the minimum possible value for the radius of a ball being tangent to $p$ disjoint unit balls in $\mathbb{R}^d$ for $d \geq 2$. Notice that $r_d^3 = r_d^*$. We can generalize the above result for appropriate $p$ as follows.

The algorithm is similar to PRECLUSTERING-DD, but in Step 2 we replace $b_{2k+1}$ by $b_{(p-1)k+1}$ and in Step 3 we return the preclustering $\mathcal{C} := \{\{b_1\}, \ldots, \{b_{(p-1)k}\}, B_1, \ldots, B_k\}$. Note that Lemmas 3 and 4 still hold if we replace $2k + 1$ with $(p - 1)k + 1$ and $r_d^*$ with $r_d^p$. Detailed proofs for the generalized lemmas are presented in the appendix.

**Theorem 7** *Let B be a set of disjoint balls in $\mathbb{R}^d$ with $d \geq 2$. Then*

*(i) there exists a $(pk, 1+\frac{2}{r_d^p})$-preclustering for the k-center and the k-median problem.*

*(ii) there exists a $(pk, \left(1 + \frac{2}{r_d^p}\right)^2)$-preclustering for the k-means problem.*

*Moreover, $(pk, 1 + \frac{2}{r_d^p} + \varepsilon)$-preclusterings for the k-center problem and for the k-median problem, can be computed in polynomial time. For the k-means problem, a $(pk, \left(1 + \frac{2}{r_d^p}\right)^2 + \varepsilon)$-preclustering can be computed in polynomial time when d is a constant.*

For instance, if we set $d = 2$ and $p = 6$, then we have $r_2^6 = 1$ since the following observation is well known. (For other bounds on $r_d^p$, see at [23].)

**Observation 8** *In $\mathbb{R}^2$, for any six disjoint unit balls, the radius of any ball intersecting all the six balls, is at least 1.*

Indeed, the familiar configuration in which a unit ball touches six unit balls around it, shows that $r_2^6$ is not greater than one. Theorem 7 together with Observation 8 lead us to the following corollary.

**Corollary 1** *Any set of disjoint balls in $\mathbb{R}^2$ admits a $(6k, 3)$-preclustering for the k-center and the k-median problem, and a $(6k, 9)$-preclustering for k-means problem.*

## 4 Approximation by an Arbitrarily Small Factor

In this section we study $(f(k), \varepsilon)$-preclusterings for approximation factor $\varepsilon < 1$.

We start by proving that for an arbitrarily small $\varepsilon$, there is no $(f(k), \varepsilon)$-preclustering for the $k$-median and the $k$-means problem, where $f(k)$ is independent from the number
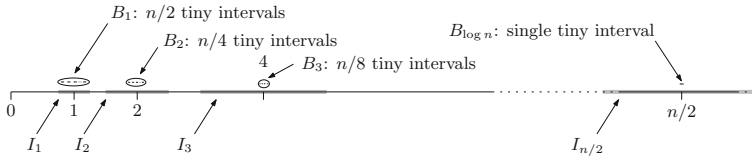
**Fig. 4** Illustration for the proof of Theorem 9. The intervals $I_i$ are indicated in grey

of imprecise points. More precisely, we prove that for an arbitrarily small $\varepsilon$, in any $(f(k), \varepsilon)$-preclustering for the $k$-median and the $k$-means problem, we have $f(k) = \Omega(\log n)$, where $n$ is the number of imprecise points.

**Theorem 9** *There exists a set $B$ of $n$ disjoint balls in $\mathbb{R}^1$ and a real number $\varepsilon > 0$ such that for the $k$-median and the $k$-means problem, in any $(f(k), \varepsilon)$-preclustering of $B$, we have $f(k) = \Omega(\log n)$.*

**Proof** We present the proof for the $k$-median problem. The proof for the $k$-means problem is similar with some minor changes, and is presented in the appendix.

Observe that it suffices to prove the lower bound for $k = 1$; for larger $k$ we can simply copy the construction $k$ times and put the copies sufficiently far from each other. Now for $k = 1$, let $n = 2^t - 1$ and for every $1 \le i \le t$ let $B_i$ be a set of $2^{t-i}$ tiny intervals all very close to the point $2^{i-1}$ in $\mathbb{R}^1$. Define $B$ as the union of $B_1, B_2, \cdots, B_t$. It is not difficult to see that $\text{OPT}_1(B, 1) = (t-2) \cdot 2^{t-1} + 1$ (which is achieved e.g. by taking the point $x = 1$ as the center). On the other hand assume that we precluster $B$ into $f(1)$ preclusters, using $q_1, q_2, \cdots, q_{f(1)}$ as centers. For $1 \le i \le t$, let $I_i$ be an open interval of length $2^{i-2}$ whose midpoint is $2^{i-1}$. Note that $I_1, I_2, \cdots I_t$ are disjoint. Thus, at least $t - f(1)$ of the intervals $I_1, I_2, \cdots I_t$ do not contain any of the centers $q_i$. But when $I_i$ does not contain a center, this means the total cost for $B_i$ is at least $2^{t-i} \cdot 2^{i-3} = 2^{t-3}$. Therefore, the total cost for any preclustering of $B$ into $f(1)$ preclusters is at least $(t - f(1)) \cdot 2^{t-3}$. For an $(f(1), \varepsilon)$-preclustering we should have

$$(t - f(1)) \cdot 2^{t-3} \le \varepsilon \cdot ((t-2) \cdot 2^{t-1} + 1).$$

This shows that $f(1) = \Omega(t)$ since otherwise we can set $\varepsilon$ to be small enough to contradict the above inequality. $\qquad\square$
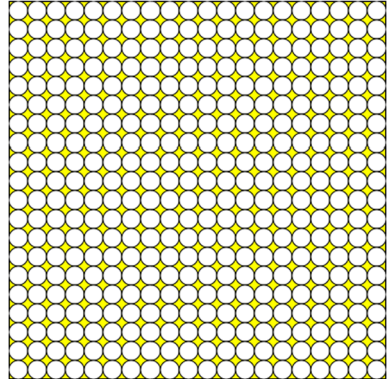
Theorem 9 shows that there is no chance to approximate the $k$-median and the $k$-means problem by any $\varepsilon$ using a constant number of preclusters. However, fortunately this is not the case for the $k$-center problem. Next, we explain how to obtain a $(\lceil \sqrt{d}/\varepsilon \rceil^d \cdot k, \varepsilon)$-preclustering for the $k$-center problem, by adding more steps to the algorithm PRECLUSTERING-DD$(B, k)$. We need the following lemma.

**Lemma 10** *For any point set $P$ in $\mathbb{R}^d$, any integer $k \ge 1$, and any $\varepsilon > 0$ we have*

$$\text{OPT}_\infty(P, c_d(\varepsilon) \cdot k) \le \varepsilon \cdot \text{OPT}_\infty(P, k)$$

*for $c_d(\varepsilon) = \lceil \sqrt{d}/\varepsilon \rceil^d$.*

**Fig. 5** $n$ unit balls forming a square in 2D



**Proof** First consider the case $k = 1$. Let $q$ be the optimal centroid for $P$ and let $S$ be the smallest hypercube centered at $q$ and containing $P$. Note that the edge length of $S$ is at most $2 \cdot \text{OPT}_\infty(P, 1)$. Partition $S$ into $\lceil \sqrt{d}/\varepsilon \rceil^d$ smaller hypercubes of edge length at most $2\varepsilon \cdot \text{OPT}_\infty(P, 1)/\sqrt{d}$, and for each such hypercube make a cluster containing all points in it. Note that each such cluster can be covered by a ball of radius $\varepsilon \cdot \text{OPT}_\infty(P, 1)$. Hence,

$$\text{OPT}_\infty(P, \lceil \sqrt{d}/\varepsilon \rceil^d) \leq \varepsilon \cdot \text{OPT}_\infty(P, 1).$$

For $k > 1$ we can simply apply the result for $k = 1$ to each of the $k$ clusters in an optimal $k$-center clustering on $P$. □

With this lemma in hand we can now run algorithm PRECLUSTERING-DD$(B, k')$ for $d \geq 2$ with the appropriate value of $k'$, namely $k' = c_d(\varepsilon/(7 + 4\sqrt{3})) \cdot k$ and then by Theorem 6 we get a $(3k', \varepsilon)$-preclustering with $k' = \Theta(\lceil \sqrt{d}/\varepsilon \rceil^d \cdot k)$. Also for $d = 1$ we can run algorithm PRECLUSTERING-1D$(B, k')$ with $k' = c_d(\varepsilon) \cdot k$ and then by Theorem 2 we get a $(3k', \varepsilon)$-preclustering with $k' = \lceil \sqrt{d}/\varepsilon \rceil^d \cdot k$.

**Theorem 11** *For any set $B$ of disjoint balls in $\mathbb{R}^d$ there exists a $(\Theta(\lceil \sqrt{d}/\varepsilon \rceil^d \cdot k), \varepsilon)$-preclustering for any positive constant $\varepsilon$.*
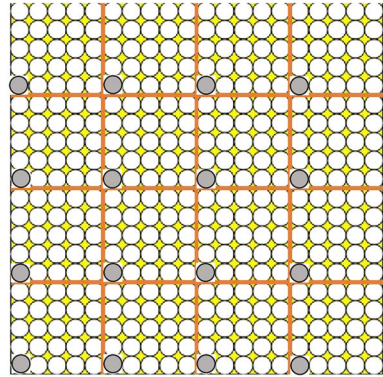
Finally, we show that this number of preclusters is in range of $O(d^d)$ from best possible cluster and as we consider our problem in fixed diamension, its asymptotically very close to the best number we can achieve for the $k$-center problem.

**Theorem 12** *There exists a set $B$ of $n$ disjoint balls in $\mathbb{R}^d$ such that in any $(f(k), \varepsilon)$-preclustering of $B$ for the $k$-center problem, we have $f(k) = \Omega(\lceil 1/(\varepsilon\sqrt{d}) \rceil^d \cdot k)$.*

**Proof** Observe that it suffices to prove the lower bound for $k = 1$; for larger $k$ we can simply copy the construction $k$ times and put the copies sufficiently far from each other. Now, for $k = 1$ consider a set $B$ of $n^{1/d} \times \cdots \times n^{1/d}$ unit balls arranged in a grid-like pattern, as in Fig. 5. Note that

$$\text{OPT}_\infty(P, 1) \leq \sqrt{d} \cdot n^{1/d}$$

**Fig. 6** Clustering grid circles
into square-shaped clusters



for any $B$-instance $P$. Now partition the "grid" into $\lceil 1/(\varepsilon\sqrt{d})\rceil^d$ "subgrids" as in Fig. 6. For each subgrid, select the ball with the lexicographically smallest center (shaded in Fig. 6), and let $B^* \subset B$ be the set of selected balls. If a preclustering uses fewer than $\lceil 1/(\varepsilon\sqrt{d})\rceil^d$ preclusters, two of the balls from $B^*$ will end up in the same precluster. But then there is a $B$-instance $P$ where $\mathcal{C}\text{-}\mathrm{COST}_\infty(P) > \varepsilon \cdot \sqrt{d} \cdot n^{1/d}$. Hence, any $(f(1), \varepsilon)$-preclustering must have $\Omega(\lceil 1/(\varepsilon\sqrt{d})\rceil^d)$ preclusters. $\qquad\square$

## 5 Concluding Remarks

In this paper, we introduced the concept of preclustering for imprecise points and studied it for the $k$-center, $k$-median and $k$-means problems. It would be interesting if one can fill the gap between lower and upper bounds for the number of preclusters needed in order to approximate the optimum solution. Also studying the problem for not necessarily disjoint balls would be interesting. Note that in this case one has to restrict the amount of overlap by a parameter, in order to find non-trivial results.

## A Generalizations of Lemmas 3 and 4

We present the generalizations of Lemmas 3 and 4 which are used in order to prove Theorem 7. Note that $p$ can take any value between 1 and $n/k$.

**Lemma 13** *For any $B$-instance $P$ the preclustering $\mathcal{C} := \{\{b_1\}, \dots, \{b_{(p-1)k}\}, B_1, \dots, B_k\}$ computed by the generalized algorithm satisfies:*

*(i)* $\mathcal{C}\text{-}\mathrm{COST}_\infty(P) \leq \mathrm{OPT}_\infty(P, k) + 2 \cdot \mathrm{radius}(b_{(p-1)k+1})$
*(ii)* $\mathcal{C}\text{-}\mathrm{COST}_1(P) \leq \mathrm{OPT}_1(P, k) + 2\sum_{j=(p-1)k+1}^{n} \mathrm{radius}(b_j)$
*(iii)* $\sqrt{\mathcal{C}\text{-}\mathrm{COST}_2(P)} \leq \sqrt{\mathrm{OPT}_2(P, k)} + 2\sqrt{\sum_{j=(p-1)k+1}^{n} \mathrm{radius}(b_j)^2}.$

***Proof*** We first prove part (i) of the lemma. Let $P$ be any $B$-instance, let $p_j \in P$ denote the point inside $b_j$, and let $c_j$ be the center of $b_j$. Recall that $P_i \subset P$ is the subset of points in the instance corresponding to the precluster $B_i$. Define $P_{\mathrm{small}} :=$

$\{p_{(p-1)k+1}, \ldots, p_n\}$ to be the set of points from $P$ in the small balls, and define $C_{\text{small}} := \{c_{(p-1)k+1}, \ldots, c_n\}$. Note that $P_{\text{small}} = P_1 \cup \cdots \cup P_k$ and that

$$|p_j c_j| \leq \text{radius}(b_j) \leq \text{radius}(b_{(p-1)k+1}) \tag{7}$$

for all $p_j \in P_{\text{small}}$. We define the following sets of centroids:

– Let $Q := \{q_1, \ldots, q_k\}$ be the set of centroids in an optimal $k$-center solution for the entire point set $P$. We have

$$\max_{p_j \in P_{\text{small}}} \min_{q_i \in Q} |p_j q_i| \leq \max_{p_j \in P} \min_{q_i \in Q} |p_j q_i| = \text{OPT}_\infty(P, k). \tag{8}$$

– Let $Q' := \{q'_1, \ldots, q'_k\}$ be the set of centroids in the optimal $k$-center clustering on $C_{\text{small}}$ used in Step 2 of the algorithm. Thus

$$\max_{c_i \in C_{\text{small}}} \min_{q'_j \in Q'} |c_i q'_j| = \text{OPT}_\infty(C_{\text{small}}, k) \leq \max_{c_i \in C_{\text{small}}} \min_{q_j \in Q} |c_i q_j|. \tag{9}$$

– Let $Q'' := \{q''_1, \ldots, q''_k\}$, where $q''_i$ is the optimal centroid for $P_i$. Note that for all $P_i$ we have

$$\max_{p_j \in P_i} |p_j q''_i| \leq \max_{p_j \in P_i} |p_j q'_i|. \tag{10}$$

Since the total cost of the singleton preclusters is trivially zero, we have

$$
\begin{aligned}
&\mathcal{C}\text{-}\text{COST}_\infty(P) \\
&= \max_{1 \leq i \leq k} \max_{p_j \in P_i} |p_j q''_i| \\
&\leq \max_{1 \leq i \leq k} \max_{p_j \in P_i} |p_j q'_i| && \text{(Inequality (10))} \\
&\leq \max_{1 \leq i \leq k} \max_{p_j \in P_i} \left( |p_j c_j| + |c_j q'_i| \right) && \text{(triangle inequality)} \\
&\leq \text{radius}(b_{(p-1)k+1}) + \max_{1 \leq i \leq k} \max_{p_j \in P_i} |c_j q'_i| && \text{(Inequality (7))} \\
&\leq \text{radius}(b_{(p-1)k+1}) + \max_{c_j \in C_{\text{small}}} \min_{q'_i \in Q'} |c_j q'_i| && \text{(definition of } C_{\text{small}}) \\
&\leq \text{radius}(b_{(p-1)k+1}) + \max_{c_j \in C_{\text{small}}} \min_{q_i \in Q} |c_j q_i| && \text{(Inequality (9))} \\
&\leq \text{radius}(b_{(p-1)k+1}) + \max_{p_j \in P_{\text{small}}} \min_{q_i \in Q} \left( |c_j p_j| + |p_j q_i| \right) && \text{(triangle inequality)} \\
&\leq 2 \cdot \text{radius}(b_{(p-1)k+1}) + \max_{p_j \in P_{\text{small}}} \min_{q_i \in Q} |p_j q_i| && \text{(Inequality (7))} \\
&\leq 2 \cdot \text{radius}(b_{(p-1)k+1}) + \text{OPT}_\infty(P, k) && \text{(Inequality (8))}
\end{aligned}
$$

To prove part (ii) of the lemma, which deals with the $k$-median problem, note that Inequality (2) still holds while Inequalities (3)–(5) hold if we replace the max-operator by a summation. Part (ii) can thus be derived using a similar derivation as for part (i).

To prove part (iii), which deals with the $k$-means problem, we need to work with squared distances. Note that Inequality (2) still holds, while Inequalities (3)–(5) hold if we replace the max-operator with a summation and all distance values with their squared values. For squared distances the triangle inequality does not hold. Instead we use the triangle inequality for $L_2$ norm, which is called Minkowsky Inequality. A

similar computation as above can now be used to prove part (iii); we have

$$\sqrt{\mathcal{C}\text{-Cost}_2(P)}$$
$$= \sqrt{\sum_{i=1}^{k}\sum_{p_j \in P_i} |p_j q_i''|^2}$$
$$\leq \sqrt{\sum_{i=1}^{k}\sum_{p_j \in P_i} |p_j q_i'|^2} \qquad\qquad\qquad\qquad \text{(Inequality (10))}$$
$$\leq \sqrt{\sum_{i=1}^{k}\sum_{p_j \in P_i} \left(|p_j c_j| + |c_j q_i'|\right)^2} \qquad\qquad \text{(triangle inequality)}$$
$$\leq \sqrt{\sum_{i=1}^{k}\sum_{p_j \in P_i} |p_j c_j|^2} + \sqrt{\sum_{i=1}^{k}\sum_{p_j \in P_i} |c_j q_i'|^2} \qquad \text{(Minkowsky inequality)}$$
$$\leq \sqrt{\sum_{j=(p-1)k+1}^{n} \text{radius}(b_j)^2} + \sqrt{\sum_{i=1}^{k}\sum_{p_j \in P_i} |c_j q_i'|^2} \quad \text{(Inequality (7))}$$

On the other hand

$$\sqrt{\sum_{i=1}^{k}\sum_{p_j \in P_i} |c_j q_i'|^2}$$
$$\leq \sqrt{\sum_{c_j \in C_{\text{small}}} \min_{q_i' \in Q'} |c_j q_i'|^2} \qquad\qquad\qquad \text{(definition of } C_{\text{small}}\text{)}$$
$$\leq \sqrt{\sum_{c_j \in C_{\text{small}}} \min_{q_i \in Q} |c_j q_i|^2} \qquad\qquad\qquad \text{(Inequality (9))}$$
$$\leq \sqrt{\sum_{c_j \in C_{\text{small}}} \min_{q_i \in Q} \left(|c_j p_j| + |p_j q_i|\right)^2} \qquad \text{(triangle inequality)}$$
$$\leq \sqrt{\sum_{c_j \in C_{\text{small}}} |c_j p_j|^2} + \sqrt{\sum_{c_j \in C_{\text{small}}} \min_{q_i \in Q} |p_j q_i|^2} \quad \text{(Minkowsky inequality)}$$
$$\leq \sqrt{\sum_{j=(p-1)k+1}^{n} \text{radius}(b_j)^2} + \sqrt{\sum_{c_j \in C_{\text{small}}} \min_{q_i \in Q} |p_j q_i|^2} \quad \text{(Inequality (7))}$$
$$\leq \sqrt{\sum_{j=(p-1)k+1}^{n} \text{radius}(b_j)^2} + \sqrt{\text{Opt}_2(P, k)} \qquad \text{(Inequality (8))}$$

By adding up the above inequalities we conclude part (iii) of the lemma.  $\square$

**Lemma 14** *Let $r_d^p$ be the smallest possible radius of any ball that intersects $p$ disjoint unit balls in $\mathbb{R}^d$. Then*

*(i)* $\text{Opt}_\infty(P, k) \geq r_d^p \cdot \text{radius}(b_{(p-1)k+1})$
*(ii)* $\text{Opt}_1(P, k) \geq r_d^p \cdot \sum_{j=(p-1)k+1}^{n} \text{radius}(b_j)$
*(iii)* $\text{Opt}_2(P, k) \geq (r_d^p)^2 \cdot \sum_{j=(p-1)k+1}^{n} \text{radius}(b_j)^2$

**Proof** For part (i) notice that by the Pigeonhole Principle an optimal clustering must have a cluster containing at least $p$ points from $\{p_1, \ldots, p_{(p-1)k+1}\}$. The cost of this cluster is lower bounded by the radius of the smallest ball intersecting $p$ balls of radius at least $b_{(p-1)k+1}$, which is in turn lower bounded by $r_d^p \cdot \text{radius}(b_{(p-1)k+1})$.

For part (ii) let $P_1, P_2, \ldots, P_k$ be the clusters in an optimal $k$-median clustering on $P$, and let $q_i$ be the centroid of $P_i$ in this clustering. Let $B_i$ be the set of balls corresponding the points in $P_i$. We claim that

$$\sum_{p_j \in P_i} |p_j q_i| \geq r_d^p \cdot \left(\left(\sum_{b_j \in B_i} \text{radius}(b_j)\right) - s_i\right). \tag{11}$$

where $s_i$ is the sum of the radii of the $p-1$ largest balls in $B_i$. To show this, let $b(q_i, r)$ be the ball of radius $r$ centered at $q_i$, let $P_i(r) := \{p_j \in P_i : b_j \cap b(q_i, r) \neq \emptyset\}$ be

the set of points in $P_i$ whose associated ball intersects $b(q_i, r)$, and let $B_i(r)$ be the corresponding set of balls. Since for sufficiently large $r$ we have $P_i = P_i(r)$, it suffices to show that for all $r > 0$ we have

$$\sum_{p_j \in P_i(r)} |p_j q_i| \geq r_d^p \cdot \left( \left( \sum_{p_j \in B_i(r)} \text{radius}(b_j) \right) - s_i(r) \right).$$

where $s_i(r)$ is the sum of the radii of the $p - 1$ largest balls in $B_i(r)$. To prove this, consider this inequality as $r$ increases from $r = 0$ to $r = \infty$. As long as $|P_i(0)| \leq 2$ the right-hand side is zero and so the inequality is obviously true. As we increase $r$ further, $b(q_i, r)$ starts intersecting more and more balls from $B_i$. Consider what happens to the inequality when $b(q_i, r)$ starts intersecting another ball $b_\ell \in B_i$. Then $p_\ell$ is added to $P_i(r)$, so the left-hand side of the inequality increases by $|p_\ell q_i|$, which is at least $r$. The right-hand side increases by at most $r_d^p$ times the radius of the $p$-th largest ball in $B_i(r)$. By definition of $r_d^p$, if $p$ balls intersect a ball of radius $r$ then the smallest has radius at most $r/r_d^p$. Hence, the right-hand side increases by at most $r$ and the inequality remains true.

Recall that $b_1, \ldots, b_{(p-1)k}$ are the $(p - 1)k$ largest balls in $B$. Hence, summing Inequality (11) over all clusters $P_1, \ldots, P_k$ gives

$$\begin{aligned} \text{OPT}_1(P, k) &= \sum_{i=1}^{k} \sum_{p_j \in P_i} |p_j q_i| \\ &\geq r_d^p \cdot \left( \sum_{i=1}^{k} \sum_{b_j \in B_i} \text{radius}(b_j) - \sum_{j=1}^{(p-1)k} \text{radius}(b_j) \right) \\ &= r_d^p \cdot \sum_{j=(p-1)k+1}^{n} \text{radius}(b_j). \end{aligned}$$

For part (iii) the same proof as (ii) works if we replace all distances with squared distances. $\square$

## B Proof of Theorem 9 for the k-Means Problem

For the $k$-means problem, observe that it suffices to prove the lower bound for $k = 1$; for larger $k$ we can simply copy the construction $k$ times and put the copies sufficiently far from each other. Now for $k = 1$, let $n = \frac{2^{2t+1}-2}{3}$ and for every $1 \leq i \leq t$ let $B_i$ be the set of $2^{2t-2i}$ tiny intervals all very close to the point $2^{i-1}$ in $\mathbb{R}^1$. Also let $B_{-i}$ be the set of $2^{2t-2i}$ tiny intervals all very close to the point $-2^{i-1}$. Consider $B$ as the union of $B_{-t}, \cdots, B_{-2}, B_{-1}, B_1, B_2, \cdots, B_t$. It is not difficult to see that $\text{OPT}_2(B, 1) = 2t \cdot 2^{2t-2}$ (considering the origin as the center). On the other hand assume that we precluster $B$ into $f(1)$ preclusters, using $q_1, q_2, \cdots, q_{f(1)}$ as centers. For $1 \leq i \leq t$, let $I_i$ be an open interval of length $2^{i-2}$ whose midpoint is $2^{i-1}$ and let $I_{-i}$ be an open interval of length $2^{i-2}$ whose midpoint is $-2^{i-1}$. Note that $I_{-t}, \cdots, I_{-2}, I_{-1}, I_1, I_2, \cdots I_t$ are disjoint. Thus, at least $2t - f(1)$ number of the intervals $I_{-t}, \cdots, I_{-2}, I_{-1}, I_1, I_2, \cdots I_t$ are empty of a center. But when $I_i$ (similarly $I_{-i}$) is empty of a center, this means the total cost for $B_i$ (similarly $B_{-i}$) is at least $2^{2t-2i} \cdot (2^{i-3})^2) = 2^{2t-6}$. Therefore, the total cost for any preclustering of $B$ into $f(1)$ preclusters is at least $(2t - f(1)) \cdot 2^{2t-6}$. For a $(f(1), \varepsilon)$-preclustering we should

have

$$(2t - f(1)) \cdot 2^{2t-6} \le \varepsilon \cdot 2t \cdot 2^{2t-2}$$

which concludes $f(1) = \Omega(t)$.

One can easily generalize this example to any $\mathbb{R}^d$ by considering the balls with diameters as the intervals in $B$ introduced above.

## References

1. Agarwal, P.K., Procopiuc, C.M.: Exact and approximation algorithms for clustering. Algorithmica **33**(2), 201–226 (2002). https://doi.org/10.1007/s00453-001-0110-y
2. Arora, S., Raghavan, P., Rao, S.: Approximation schemes for euclidean $k$-medians and related problems. In: Proceedings of the ACM Symposium on the Theory of Computing, pp. 106–113 (1998). https://doi.org/10.1145/276698.276718
3. Bilu, Y., Linial, N.: Are stable instances easy? Comb. Probab. Comput. **21**(5), 643–660 (2012). https://doi.org/10.1017/S0963548312000193
4. Buchin, K., Löffler, M., Morin, P., Mulzer, W.: Preprocessing imprecise points for delaunay triangulation: simplified and extended. Algorithmica **61**(3), 674–693 (2011). https://doi.org/10.1007/s00453-010-9430-0
5. Cohen-Addad, V., Schwiegelshohn, C.: On the local structure of stable clustering instances. In: C. Umans (ed.) 58th IEEE Annual Symposium on Foundations of Computer Science, FOCS 2017, Berkeley, CA, USA, October 15–17, 2017, pp. 49–60. IEEE Computer Society (2017). https://doi.org/10.1109/FOCS.2017.14
6. Feldman, D., Monemizadeh, M., Sohler, C.: A PTAS for k-means clustering based on weak coresets. In: Proceedings of ACM Symposium on Computational Geometry, pp. 11–18 (2007). https://doi.org/10.1145/1247069.1247072
7. Gullo, F., Tagarelli, A.: Uncertain centroid based partitional clustering of uncertain data. Proc. VLDB Endow. 5(7), 610–621 (2012). https://doi.org/10.14778/2180912.2180914
8. Hochbaum, D.S., Shmoys, D.B.: A best possible heuristic for the k-center problem. Math. Oper. Res. **10**(2), 180–184 (1985). https://doi.org/10.1287/moor.10.2.180
9. Jaiswal, R., Kumar, A., Sen, S.: A simple D 2-sampling based PTAS for k-means and other clustering problems. In: Computing and Combinatorics COCOON 2012. Lecture Notes in Computer Science, vol. 7434, pp. 13–24. Springer (2012). https://doi.org/10.1007/978-3-642-32241-9_2
10. Ju, W., Luo, J., Zhu, B., Daescu, O.: Largest area convex hull of imprecise data based on axis-aligned squares. J. Comb. Optim. **26**(4), 832–859 (2013). https://doi.org/10.1007/s10878-012-9488-5
11. Kao, B., Lee, S.D., Cheung, D.W., Ho, W., Chan, K.F.: Clustering uncertain data using voronoi diagrams. In: Proceedings of the 8th IEEE International Conference on Data Mining (ICDM 2008), December 15–19, 2008, Pisa, Italy, pp. 333–342. IEEE Computer Society (2008). https://doi.org/10.1109/ICDM.2008.31
12. Kolliopoulos, S.G., Rao, S.: A nearly linear-time approximation scheme for the euclidean kappamedian problem. In: Algorithms—ESA Proceedings. Lecture Notes in Computer Science, vol. 1643, pp. 378–389. Springer (1999). https://doi.org/10.1007/3-540-48481-7_33
13. Liu, C., Montanari, S.: Minimizing the diameter of a spanning tree for imprecise points. Algorithmica **80**(2), 801–826 (2018). https://doi.org/10.1007/s00453-017-0292-6
14. Löffler, M.: Data imprecision in computational geometry. Ph.D. thesis, Utrecht University, Netherlands (2009)
15. Löffler, M., van Kreveld, M.J.: Largest and smallest convex hulls for imprecise points. Algorithmica **56**(2), 235–269 (2010). https://doi.org/10.1007/s00453-008-9174-2
16. Mahajan, M., Nimbhorkar, P., Varadarajan, K.R.: The planar k-means problem is nphard. Theor. Comput. Sci. **442**, 13–21 (2012). https://doi.org/10.1016/j.tcs.2010.05.034
17. Megiddo, N., Supowit, K.J.: On the complexity of some common geometric location problems. SIAM J. Comput. **13**(1), 182–196 (1984). https://doi.org/10.1137/0213014

18. Nagai, T., Tokura, N.: Tight error bounds of geometric problems on convex objects with imprecise coordinates. In: JCDCG. Lecture Notes in Computer Science, vol. 2098, pp. 252–263. Springer (2000). https://doi.org/10.1007/3-540-47738-1_24

19. Saulpic, D., Cohen-Addad, V., Feldmann, A.E.: Near-linear time approximations schemes for clustering in doubling metrics. In: D. Zuckerman (ed.) 60th IEEE Annual Symposium on Foundations of Computer Science, FOCS 2019, Baltimore, Maryland, USA, November 9–12, 2019, pp. 540–559. IEEE Computer Society (2019). https://doi.org/10.1109/FOCS.2019.00041

20. Sheikhi, F., Mohades, A., de Berg, M., Mehrabi, A.D.: Separability of imprecise points. Comput. Geom. **61**, 24–37 (2017). https://doi.org/10.1016/j.comgeo.2016.10.001

21. Shmoys, D.B., Tardos, É.: An approximation algorithm for the generalized assignment problem. Math. Program. **62**, 461–474 (1993). https://doi.org/10.1007/BF01585178

22. Spielman, D.A., Teng, S.: Smoothed analysis: an attempt to explain the behavior of algorithms in practice. Commun. ACM **52**(10), 76–84 (2009). https://doi.org/10.1145/1562764.1562785

23. Talata, I.: Exponential lower bound for the translative kissing numbers of d-dimensional convex bodies. Discrete Comput. Geom. **19**(3), 447–455 (1998). https://doi.org/10.1007/PL00009362