# Range-Clustering Queries[*][†]

## Mikkel Abrahamsen[1], Mark de Berg[2], Kevin Buchin[3], Mehran Mehr[4], and Ali D. Mehrabi[5]

1     **Department of Computer Science, University of Copenhagen, Copenhagen, Denmark**
    `miab@di.ku.dk`
2     **Department of Computer Science, TU Eindhoven, Eindhoven, The Netherlands**
    `mdberg@win.tue.nl`
2     **Department of Computer Science, TU Eindhoven, Eindhoven, The Netherlands**
    `k.a.buchin@tue.nl`
2     **Department of Computer Science, TU Eindhoven, Eindhoven, The Netherlands**
    `m.mehr@tue.nl`
2     **Department of Computer Science, TU Eindhoven, Eindhoven, The Netherlands**
    `amehrabi@win.tue.nl`

## Abstract

In a geometric $k$-clustering problem the goal is to partition a set of points in $\mathbb{R}^d$ into $k$ subsets such that a certain cost function of the clustering is minimized. We present data structures for orthogonal *range-clustering queries* on a point set $S$: given a query box $Q$ and an integer $k \geqslant 2$, compute an optimal $k$-clustering for $S \cap Q$. We obtain the following results.

- We present a general method to compute a $(1+\varepsilon)$-approximation to a range-clustering query, where $\varepsilon > 0$ is a parameter that can be specified as part of the query. Our method applies to a large class of clustering problems, including $k$-center clustering in any $L_p$-metric and a variant of $k$-center clustering where the goal is to minimize the sum (instead of maximum) of the cluster sizes.
- We extend our method to deal with capacitated $k$-clustering problems, where each of the clusters should not contain more than a given number of points.
- For the special cases of rectilinear $k$-center clustering in $\mathbb{R}^1$, and in $\mathbb{R}^2$ for $k = 2$ or $3$, we present data structures that answer range-clustering queries exactly.

## 1   Introduction

**Motivation**

The range-searching problem is one of the most important and widely studied problems in computational geometry. In the standard setting one is given a set $S$ of points in $\mathbb{R}^d$, and a query asks to report or count all points inside a geometric query range $Q$. In many

---

applications, however, one would like to perform further analysis on the set $S \cap Q$, to obtain more information about its structure. Currently one then has to proceed as follows: first perform a range-reporting query to explicitly report $S \cap Q$, then apply a suitable analysis algorithm to $S \cap Q$. This two-stage process can be quite costly, because algorithms for analyzing geometric data sets can be slow and $S \cap Q$ can be large. To avoid this we would need data structures for what we call *range-analysis queries*, which directly compute the desired structural information about $S \cap Q$. In this paper we develop such data structures for the case where one is interested in a cluster-analysis of $S \cap Q$.

Clustering is a fundamental task in data analysis. It involves partitioning a given data set into subsets called *clusters*, such that similar elements end up in the same cluster. Often the data elements can be viewed as points in a geometric space, and similarity is measured by considering the distance between the points. We focus on clustering problems of the following type. Let $S$ be a set of $n$ points in $\mathbb{R}^d$, and let $k \geqslant 2$ be a natural number. A *k-clustering* of $S$ is a partitioning $\mathcal{C}$ of $S$ into at most $k$ clusters. Let $\Phi(\mathcal{C})$ denote the *cost* of $\mathcal{C}$. The goal is now to find a clustering $\mathcal{C}$ that minimizes $\Phi(\mathcal{C})$. Many well-known geometric clustering problems are of this type. Among them is the $k$-center problem. In the *Euclidean k-center problem* $\Phi(\mathcal{C})$ is the maximum cost of any of the clusters $C \in \mathcal{C}$, where the cost of $C$ is the radius of its smallest enclosing ball. Hence, in the Euclidean $k$-center problem we want to cover the point set $S$ by $k$ congruent balls of minimum radius. The *rectilinear k-center problem* is defined similarly except that one considers the $L_\infty$-metric; thus we want to cover $S$ by $k$ congruent axis-aligned cubes[1] of minimum size. The $k$-center problem, including the important special case of the 2-center problem, has been studied extensively, both for the Euclidean case (e.g. [2, 8, 12, 17, 16, 22]) and for the rectilinear case (e.g. [7, 23]).

All papers mentioned above – in fact, all papers on clustering that we know of – consider clustering in the single-shot version. We are the first to study *range-clustering queries* on a point set $S$: given a query range $Q$ and a parameter $k$, solve the given $k$-clustering problem on $S \cap Q$. We study this problem for the case where the query range is an axis-aligned box.

### Related work

Range-analysis queries can be seen as a very general form of range-aggregate queries. In a range-aggregate query, the goal is to compute some aggregate function $F(S \cap Q)$ over the points in the query range. The current state of the art typically deals with simple aggregate functions of the following form: each point $p \in S$ has a *weight* $w(p) \in \mathbb{R}$, and $F(S \cap Q) := \bigoplus_{p \in S \cap Q} w(p)$, where $\oplus$ is a semi-group operation. Such aggregate functions are *decomposable*, meaning that $F(A \cap B)$ can be computed from $F(A)$ and $F(B)$, which makes them easy to handle using existing data structures such as range trees.

Only some, mostly recent, papers describe data structures supporting non-decomposable analysis tasks. Several deal with finding the closest pair inside a query range (e.g. [1, 10, 13]). However, the closest pair does not give information about the global shape or distribution of $S \cap Q$, which is what our queries are about. The recent works by Brass *et al.* [5] and by Arya *et al.* [4] are more related to our paper. Brass *et al.* [5] present data structures for finding extent measures, such the width, area or perimeter of the convex hull of $S \cap Q$, or the smallest enclosing disk. (Khare *et al.* [18] improve the result on smallest-enclosing-disk queries.) These measures are strictly speaking not decomposable, but they depend only on

---

[1] Throughout the paper, when we speak of cubes (or squares, or rectangles, or boxes) we always mean axis-aligned cubes (or squares, or rectangles, or boxes).

the convex hull of $S \cap Q$ and convex hulls are decomposable. A related result is by Nekrich and Smid [20], who present a data structure that returns an $\varepsilon$-coreset inside a query range. The measure studied by Arya *et al.* [4], namely the length of the MST of $S \cap Q$, cannot be computed form the convex hull either: like our range-clustering queries, it requires more information about the structure of the point set. Thus our paper continues the direction set out by Arya *et al.*, which is to design data structures for more complicated analysis tasks on $S \cap Q$.

**Our contribution**

Our main result is a general method to answer *approximate* orthogonal range-clustering queries in $\mathbb{R}^d$. Here the query specifies (besides the query box $Q$ and the number of clusters $k$) a value $\varepsilon > 0$; the goal then is to compute a $k$-clustering $\mathcal{C}$ of $S \cap Q$ with $\Phi(\mathcal{C}) \leqslant (1 + \varepsilon) \cdot \Phi(\mathcal{C}_{\mathrm{opt}})$, where $\mathcal{C}_{\mathrm{opt}}$ is an optimal clustering for $S \cap Q$. Our method works by computing a sample $R \subseteq S \cap Q$ such that solving the problem on $R$ gives us the desired approximate solution. We show that for a large class of cost functions $\Phi$ we can find such a sample of size only $O(k(f(k)/\varepsilon)^d)$, where $f(k)$ is a function that only depends on the number of clusters. This is similar to the approach taken by Har-Peled and Mazumdar [15], who solve the (single-shot) approximate $k$-means and $k$-median problem efficiently by generating a coreset of size $O((k/\varepsilon^d) \cdot \log n)$. A key step in our method is a procedure to efficiently compute a lower bound on the value of an optimal solution within the query range. The class of clustering problems to which our method applies includes the $k$-center problem in any $L_p$-metric, variants of the $k$-center problem where we want to minimize the sum (rather than maximum) of the cluster radii, and the 2-dimensional problem where we want to minimize the maximum or sum of the perimeters of the clusters. Our technique allows us, for instance, to answer rectilinear $k$-center queries in the plane in $O((1/\varepsilon) \log n + 1/\varepsilon^2)$ for $k = 2$ or 3, in $O((1/\varepsilon) \log n + (1/\varepsilon^2)\mathrm{polylog}(1/\varepsilon))$ for $k = 4$ or 5, and in $O((k/\varepsilon) \log n + (k/\varepsilon)^{O(\sqrt{k})})$ time for $k > 3$. We also show that for the rectilinear (or Euclidean) $k$-center problem, our method can be extended to deal with the capacitated version of the problem. In the capacitated version each cluster should not contain more than $\alpha \cdot (|S \cap Q|/k)$ points, for a given $\alpha > 1$.

In the second part of the paper we turn our attention to exact solutions to range-clustering queries. Here we focus on rectilinear *k-center queries* – that is, range-clustering queries for the rectilinear $k$-center problem – in $\mathbb{R}^1$ and $\mathbb{R}^2$. We present two linear-size data structures for queries in $\mathbb{R}^1$; one has $O(k^2 \log^2 n)$ query time, the other has $O(3^k \log n)$ query time. For queries in $\mathbb{R}^2$ we present a data structure that answers 2-center queries in $O(\log n)$ time, and one that answers 3-center queries in $O(\log^2 n)$ time. Both data structures use $O(n \log^\varepsilon n)$ storage, where $\varepsilon > 0$ is an arbitrary small (but fixed) constant.

## 2 Approximate Range-Clustering Queries

In this section we present a general method to answer approximate range-clustering queries. We start by defining the class of clustering problems to which it applies.

Let $S$ be a set of $n$ points in $\mathbb{R}^d$ and let $\mathrm{Part}(S)$ be the set of all partitions of $S$. Let $\mathrm{Part}_k(S)$ be the set of all partitions into at most $k$ subsets, that is, all $k$-clusterings of $S$. Let $\Phi : \mathrm{Part}(S) \mapsto \mathbb{R}_{\geqslant 0}$ be the cost function defining our clustering problem, and define

$$\mathrm{OPT}_k(S) := \min_{\mathcal{C} \in \mathrm{Part}_k(S)} \Phi(\mathcal{C})$$

to be the minimum cost of any $k$-clustering. Thus the goal of a range-clustering query with query range $Q$ and parameter $k \geqslant 2$ is to compute a clustering $\mathcal{C} \in \mathrm{Part}_k(S_Q)$ such that

$\Phi(\mathcal{C}) = \text{OPT}_k(S_Q)$, where $S_Q := S \cap Q$. From now on we use $S_Q$ as a shorthand for $S \cap Q$. The method presented in this section gives an approximate answer to such a query: for a given constant $\varepsilon > 0$, which can be specified as part of the query, the method will report a clustering $\mathcal{C} \in \text{Part}_k(S_Q)$ with $\Phi(\mathcal{C}) \leqslant (1 + \varepsilon) \cdot \text{OPT}_k(S_Q)$.

To define the class of clusterings to which our method applies, we will need the concept of $r$-packings [14]. Actually, we will use a slightly weaker variant, which we define as follows. Let $|pq|$ denote the Euclidean distance between two points $p$ and $q$. A subset $R \subseteq P$ of a point set $P$ is called a *weak r-packing* for $P$, for some $r > 0$, if for any point $p \in P$ there exists a *packing point* $q \in R$ such that $|pq| \leqslant r$. (The difference with standard $r$-packings is that we do not require that $|qq'| > r$ for any two points $q, q' \in R$.) The clustering problems to which our method applies are the ones whose cost function is *regular*, as defined next.

▶ **Definition 1.** A cost function $\Phi : \text{Part}(S) \mapsto \mathbb{R}_{\geqslant 0}$ is called $(c, f(k))$-*regular*, if there is a constant $c$ and a function $f : \mathbb{N}_{\geqslant 2} \mapsto \mathbb{R}_{\geqslant 0}$ such that the followings hold.

- For any clustering $\mathcal{C} \in \text{Part}(S)$, we have

$$\Phi(\mathcal{C}) \geqslant c \cdot \max_{C \in \mathcal{C}} \text{diam}(C),$$

  where $\text{diam}(C) = \max_{p,q \in C} |pq|$ denotes the Euclidean diameter of the cluster $C$. We call this the *diameter-sensitivity* property.

- For any subset $S' \subseteq S$, any weak $r$-packing $R$ of $S'$, and any $k \geqslant 2$, we have that

$$\text{OPT}_k(R) \leqslant \text{OPT}_k(S') \leqslant \text{OPT}_k(R) + r \cdot f(k).$$

  Moreover, given a $k$-clustering $\mathcal{C} \in \text{Part}_k(R)$ we can compute a $k$-clustering $\mathcal{C}^* \in \text{Part}_k(S')$ with $\Phi(\mathcal{C}^*) \leqslant \Phi(\mathcal{C}) + r \cdot f(k)$ in time $T_{\text{expand}}(n, k)$. We call this the *expansion* property.

**Examples**

Many clustering problems have regular cost functions, in particular when the cost function is the aggregation – the sum, for instance, or the maximum – of the costs of the individual clusters. Next we give some examples.

**Rectilinear and other $k$-center problems.** For a cluster $C$, let $\text{radius}_p(C)$ denote the radius of the minimum enclosing ball of $C$ in the $L_p$-metric. In the $L_\infty$-metric, for instance, $\text{radius}_p(C)$ is half the edge length of a minimum enclosing axis-aligned cube of $C$. Then the cost of a clustering $\mathcal{C}$ for the $k$-center problem in the $L_p$-metric is $\Phi_p^{\max}(\mathcal{C}) = \max_{C \in \mathcal{C}} \text{radius}_p(C)$. One can easily verify that the cost function for the rectilinear $k$-center problem is $(1/(2\sqrt{d}), 1)$-regular, and for the Euclidean $k$-center problem it is $(1/2, 1)$-regular. Moreover, $T_{\text{expand}}(n, k) = O(k)$ for the $k$-center problem, since we just have to scale each ball by adding $r$ to its radius.[2] (In fact $\Phi_p^{\max}(\mathcal{C})$ is regular for any $p$.)

**Min-sum variants of the $k$-center problem.** In the $k$-center problem the goal is to minimize $\max_{C \in \mathcal{C}} \text{radius}_p(C)$. Instead we can also minimize $\Phi_p^{\text{sum}}(\mathcal{C}) := \sum_{C \in \mathcal{C}} \text{radius}_p(C)$, the sum of the cluster radii. Also these costs functions are regular; the only difference is that the expansion property is now satisfied with $f(k) = k$, instead of with $f(k) = 1$. Another interesting variant is to minimize $\left(\sum_{C \in \mathcal{C}} \text{radius}_2(C)^2\right)^{1/2}$, which is $(1/(2\sqrt{d}), \sqrt{k})$-regular.

---

[2] This time bound only accounts for reporting the set of balls that define the clustering. If we want to report the clusters explicitly, we need to add an $O(nk)$ term.

**Minimum perimeter $k$-clustering problems.**    For a cluster $C$ of points in $\mathbb{R}^2$, define $\mathrm{per}(C)$ to be the length of the perimeter of the convex hull of $C$. In the minimum perimeter-sum clustering problem the goal is to compute a $k$-clustering $\mathcal{C}$ such that $\Phi_{\mathrm{per}} := \sum_{C \in \mathcal{C}} \mathrm{per}(C)$ is minimized [6]. This cost function is $(2, 2\pi k)$-regular. Indeed, if we expand the polygons in a clustering $\mathcal{C}$ of a weak $r$-packing $R$ by taking the Minkowski sum with a disk of radius $r$, then the resulting shapes cover all the points in $S$. Each perimeter increases by $2\pi r$ in this process. To obtain a clustering, we then assign each point to the cluster of its closest packing point, so $T_{\mathrm{expand}}(n, k) = O(n \log n)$.

**Non-regular costs functions.**    Even though many clustering problems have regular costs functions, not all clustering problems do. For example, the $k$-means problem does not have a regular cost function. Minimizing the the max or sum of the areas of the convex hulls of the clusters is not regular either.

### Our data structure and query algorithm

We start with a high-level overview of our approach. Let $S$ be the given point set on which we want to answer range-clustering queries, and let $Q$ be the query range. We assume we have an algorithm $\textsc{SingleShotClustering}(P, k)$ that computes an optimal solution to the $k$-clustering problem (for the given cost function $\Phi$) on a given point set $P$. (Actually, it is good enough if $\textsc{SingleShotClustering}(P, k)$ gives a $(1 + \varepsilon)$-approximation.) Our query algorithm proceeds as follows.

>    $\textsc{ClusterQuery}(k, Q, \varepsilon)$
> 1:   Compute a lower bound $\textsc{lb}$ on $\textsc{Opt}_k(S_Q)$.
> 2:   Set $r := \varepsilon \cdot \textsc{lb}/f(k)$ and compute a weak $r$-packing $R_Q$ on $S_Q$.
> 3:   $\mathcal{C} := \textsc{SingleShotClustering}(R_Q, k)$.
> 4:   Expand $\mathcal{C}$ into a $k$-clustering $\mathcal{C}^*$ of cost at most $\Phi(\mathcal{C}) + r \cdot f(k)$ for $S_Q$.
> 5:   **return** $\mathcal{C}^*$.

Note that Step 4 is possible because $\Phi$ is $(c, f(k))$-regular. The following lemma is immediate.

▶ **Lemma 2.** $\Phi(\mathcal{C}^*) \leqslant (1 + \varepsilon) \cdot \textsc{Opt}_k(S_Q)$.

Next we show how to perform Steps 1 and 2: we will describe a data structure that allows us to compute a suitable lower bound $\textsc{lb}$ and a corresponding weak $r$-packing, such that the size of the $r$-packing depends only on $\varepsilon$ and $k$ but not on $|S_Q|$.

Our lower bound and $r$-packing computations are based on so-called cube covers. A *cube cover* of $S_Q$ is a collection $\mathcal{B}$ of interior-disjoint cubes that together cover all the points in $S_Q$ and such that each $B \in \mathcal{B}$ contains at least one point from $S_Q$ (in its interior or on its boundary). Define the size of a cube $B$, denoted by $\mathrm{size}(B)$, to be its edge length. The following lemma follows immediately from the fact that the diameter of a cube $B$ in $\mathbb{R}^d$ is $\sqrt{d} \cdot \mathrm{size}(B)$.

▶ **Lemma 3.** *Let $\mathcal{B}$ be a cube cover of $S_Q$ such that $\mathrm{size}(B) \leqslant r/\sqrt{d}$ for all $B \in \mathcal{B}$. Then any subset $R \subseteq S_Q$ containing a point from each cube $B \in \mathcal{B}$ is a weak $r$-packing for $S$.*

Our next lemma shows we can find a lower bound on $\textsc{Opt}_k(S_Q)$ from a suitable cube cover.

▶ **Lemma 4.** *Suppose the cost function $\Phi$ is $(c, f(k))$-regular. Let $\mathcal{B}$ be a cube cover of $S_Q$ such that $|\mathcal{B}| > k2^d$. Then $\textsc{Opt}_k(S_Q) \geqslant c \cdot \min_{B \in \mathcal{B}} \mathrm{size}(B)$.*

**Proof.** For two cubes $B, B'$ such that the maximum $x_i$-coordinate of $B$ is at most the minimum $x_i$-coordinate of $B'$, we say that $B$ is *i-below* $B'$ and $B'$ is *i-above* $B$. We denote this relation by $B \prec_i B'$. Now consider an optimal $k$-clustering $\mathcal{C}_{\mathrm{opt}}$ of $S_Q$. By the pigeonhole principle, there is a cluster $C \in \mathcal{C}_{\mathrm{opt}}$ containing points from at least $2^d + 1$ cubes. Let $\mathcal{B}_C$ be the set of cubes that contain at least one point in $C$.

Clearly, if there are cubes $B, B', B'' \in \mathcal{B}_C$ such that $B' \prec_i B \prec_i B''$ for some $1 \leqslant i \leqslant d$, then the cluster $C$ contains two points (namely from $B'$ and $B''$) at distance at least $\mathrm{size}(B)$ from each other. Since $\Phi$ is $(c, f(k))$-regular this implies that $\Phi(\mathcal{C}_{\mathrm{opt}}) \geqslant c \cdot \mathrm{size}(B)$, which proves the lemma.

Now suppose for a contradiction that such a triple $B', B, B''$ does not exist. Then we can define a characteristic vector $\Gamma(B) = (\Gamma_1(B), \ldots, \Gamma_d(B))$ for each cube $B \in \mathcal{B}_C$, as follows:

$$\Gamma_i(B) = \begin{cases} 0 & \text{if no cube } B' \in \mathcal{B}_C \text{ is } i\text{-below } B \\ 1 & \text{otherwise} \end{cases}$$

Since the number of distinct characteristic vectors is $2^d < |B_C|$, there must be two cubes $B_1, B_2 \in B_C$ with identical characteristic vectors. However, any two interior-disjoint cubes can be separated by an axis-parallel hyperplane, so there is at least one $i \in \{1, \ldots, d\}$ such that $B_1 \prec_i B_2$ or $B_2 \prec_i B_1$. Assume w.l.o.g. that $B_1 \prec_i B_2$, so $\Gamma_i(B_2) = 1$. Since $\Gamma(B_1) = \Gamma(B_2)$ there must be a cube $B_3$ with $B_3 \prec_i B_1$. But then we have a triple $B_3 \prec_i B_1 \prec_i B_2$, which is a contradiction.                    ◀

Next we show how to efficiently perform Steps 1 and 2 of CLUSTERQUERY. Our algorithm uses a compressed octree $\mathcal{T}(S)$ on the point set $S$, which we briefly describe next.

For an integer $s$, let $G_s$ denote the grid in $\mathbb{R}^d$ whose cells have size $2^s$ and for which the origin $O$ is a grid point. A *canonical cube* is any cube that is a cell of a grid $G_s$, for some integer $s$. A *compressed octree* on a point set $S$ in $\mathbb{R}^d$ contained in a canonical cube $B$ is a tree-like structure defined recursively, as follows.

- If $|S| \leqslant 1$, then $\mathcal{T}(S)$ consists of a single leaf node, which corresponds to the cube $B$.
- If $|S| > 1$, then consider the cubes $B_1, \ldots, B_{2^d}$ that result from cutting $B$ into $2^d$ equal-sized cubes.
  - If at least two of the cubes $B_i$ contain at least one point from $S$ then $\mathcal{T}(S)$ consists of a root node with $2^d$ children $v_1, \ldots, v_{2^d}$, where $v_i$ is the root of a compressed octree for[3] $B_i \cap S$.
  - If all points from $S$ lie in the same cube $B_i$, then let $B_{\mathrm{in}} \subseteq B_i$ be the smallest canonical cube containing all points in $S$. Now $\mathcal{T}(S)$ consists of a root node with two children: one child $v$ which is the root of a compressed octree for $S$ inside $B_{\mathrm{in}}$, and one leaf node $w$ which represents the donut region $B \setminus B_{\mathrm{in}}$.

A compressed octree for a set $S$ of $n$ points in any fixed dimension can be computed in $O(n \log n)$ time, assuming a model of computation where the smallest canonical cube of two points can be computed in $O(1)$ time [14, Theorem 2.23]. For a node $v \in \mathcal{T}(S)$, we denote the cube or donut corresponding to $v$ by $B_v$, and we define $S_v := B_v \cap S$. It will be convenient to slightly modify the compressed quadtree by removing all nodes $v$ such that $S_v = \emptyset$. (These nodes must be leaves.) Note that this removes all nodes $v$ such that $B_v$ is a donut. As a result, the parent of a donut node now has only one child. The modified tree $\mathcal{T}(S)$ – with a slight abuse of terminology we still refer to $\mathcal{T}(S)$ as a compressed octree

---

[3] Here we assume that points on the boundary between cubes are assigned to one of these cubes in a consistent manner.

---

**Algorithm 1** Algorithm for Steps 1 and 2 of CLUSTERQUERY, for a $(c, f(k))$-regular cost function.

---

1: $\mathcal{B}_{\text{inner}} := B_{\text{root}(\mathcal{T}(S))}$ and $\mathcal{B}_{\text{leaf}} := \emptyset$.
2: ▷ Phase 1: Compute a lower bound on $\text{OPT}_k(S_Q)$.
3: **while** $|\mathcal{B}_{\text{inner}} \cup \mathcal{B}_{\text{leaf}}| \leqslant k2^{2d}$ and $\mathcal{B}_{\text{inner}} \neq \emptyset$ **do**
4:     Remove a largest cube $B_v$ from $\mathcal{B}_{\text{inner}}$. Let $v$ be the corresponding node.
5:     **if** $B_v \nsubseteq Q$ **then**
6:         Compute $\text{bb}(S_Q \cap B_v)$, the bounding box of $S_Q \cap B_v$.
7:         Find the deepest node $u$ such that $\text{bb}(S_Q \cap B_v) \subseteq B_u$ and set $v := u$.
8:     **end if**
9:     For each child $w$ of $v$ such that $B_w \cap S_Q \neq \emptyset$, insert $B_w$ into $\mathcal{B}_{\text{inner}}$ if $w$ is an
        internal node and insert $B_w$ into $\mathcal{B}_{\text{leaf}}$ if $w$ is a leaf node.
10: **end while**
11: $\text{LB} := c \cdot \max_{B_v \in \mathcal{B}_{\text{inner}}} \text{size}(B_v)$ .
12: ▷ Phase 2: Compute a suitable weak $r$-packing.
13: $r := \varepsilon \cdot \text{LB}/f(k)$.
14: **while** $\mathcal{B}_{\text{inner}} \neq \emptyset$ **do**
15:     Remove a cube $B_v$ from $\mathcal{B}_{\text{inner}}$ and handle it as in lines 5–9, with the following
        change: if $\text{size}(B_w) \leqslant r/\sqrt{d}$ then always insert $B_w$ into $\mathcal{B}_{\text{leaf}}$ (not into $\mathcal{B}_{\text{inner}}$).
16: **end while**
17: For each cube $B_v \in \mathcal{B}_{\text{leaf}}$ pick a point in $S_Q \cap B_v$ and put it into $R_Q$.
18: **return** $R_Q$.

---

– has the property that any internal node has at least two children. We augment $\mathcal{T}(S)$ by storing at each node $v$ an arbitrary point $p \in B_v \cap S$.

Our algorithm descends into $\mathcal{T}(S)$ to find a cube cover $\mathcal{B}$ of $S_Q$ consisting of canonical cubes, such that $\mathcal{B}$ gives us a lower bound on $\text{OPT}_k(S_Q)$. In a second phase, the algorithm then refines the cubes in the cover until they are small enough so that, if we select one point from each cube, we get a weak $r$-packing of $S_Q$ for the appropriate value of $r$. The details are described in Algorithm 1, where we assume for simplicity that $|S_Q| > 1$. (The case $|S_Q| \leqslant 1$ is easy to check and handle.)

Note that we continue the loop in lines 3–9 until we collect $k2^{2d}$ cubes (and not $k2^d$, as Lemma 4 would suggest) and that in line 11 we take the maximum cube size (instead of the minimum, as Lemma 4 would suggest).

▶ **Lemma 5.** *The value* LB *computed by Algorithm 1 is a correct lower bound on* $\text{OPT}_k(S_Q)$, *and the set* $R_Q$ *is a weak $r$-packing for* $r = \varepsilon \cdot \text{LB}/f(k)$ *of size* $O(k(f(k)/(c\,\varepsilon))^d)$.

**Proof.** As the first step to prove that LB is a correct lower bound, we claim that the loop in lines 3–9 maintains the following invariant: (i) $\bigcup(\mathcal{B}_{\text{inner}} \cup \mathcal{B}_{\text{leaf}})$ contains all points in $S_Q$, and (ii) each $B \in \mathcal{B}_{\text{inner}}$ contains at least two points from $S_Q$ and each $B \in \mathcal{B}_{\text{leaf}}$ contains exactly one point from $S_Q$. This is trivially true before the loop starts, under our assumption that $|S_Q| \geqslant 2$. Now suppose we handle a cube $B_v \in \mathcal{B}_{\text{inner}}$. If $B_v \subseteq Q$ then we insert the cubes $B_w$ of all children into $\mathcal{B}_{\text{inner}}$ or $\mathcal{B}_{\text{leaf}}$, which restores the invariant. If $B_v \nsubseteq Q$ then we first replace $v$ by $u$. The condition $\text{bb}(S_Q \cap B_v) \subseteq B_u$ guarantees that all points of $S_Q$ in $B_v$ are also in $B_u$. Hence, if we then insert the cubes $B_w$ of $u$'s children into $\mathcal{B}_{\text{inner}}$ or $\mathcal{B}_{\text{leaf}}$, we restore the invariant. Thus at any time, and in particular after the loop, the set $\mathcal{B}_{\text{inner}} \cup \mathcal{B}_{\text{leaf}}$ is a cube cover of $S_Q$.

To complete the proof that LB is a correct lower bound we do not work with the set $\mathcal{B}_{\text{inner}} \cup \mathcal{B}_{\text{leaf}}$ directly, but we work with a set $\mathcal{B}$ defined as follows. For a cube $B_v \in \mathcal{B}_{\text{inner}} \cup \mathcal{B}_{\text{leaf}}$, define parent($B_v$) to be the cube $B_u$ corresponding to the parent node $u$ of $v$. For each cube $B_v \in \mathcal{B}_{\text{inner}} \cup \mathcal{B}_{\text{leaf}}$ we put one cube into $\mathcal{B}$, as follows. If there is another cube $B_w \in \mathcal{B}_{\text{inner}} \cup \mathcal{B}_{\text{leaf}}$ such that parent($B_w$) $\subsetneq$ parent($B_v$), then we put $B_v$ itself into $\mathcal{B}$, and otherwise we put parent($B_v$) into $\mathcal{B}$. Finally, we remove all duplicates from $\mathcal{B}$. Since $\mathcal{B}_{\text{inner}} \cup \mathcal{B}_{\text{leaf}}$ is a cube cover for $S_Q$ – that is, the cubes in $\mathcal{B}_{\text{inner}} \cup \mathcal{B}_{\text{leaf}}$ are disjoint and they cover all points in $S_Q$ – the same is true for $\mathcal{B}$. Moreover, the only duplicates in $\mathcal{B}$ are cubes that are the parent of multiple nodes in $\mathcal{B}_{\text{inner}} \cup \mathcal{B}_{\text{leaf}}$, and so $|\mathcal{B}| \geqslant |\mathcal{B}_{\text{inner}} \cup \mathcal{B}_{\text{leaf}}|/2^d > k2^d$. By Lemma 4 we have $\text{OPT}_k(S_Q) \geqslant c \cdot \min_{B_v \in \mathcal{B}} \text{size}(B_v)$.

It remains to argue that $\min_{B_v \in \mathcal{B}} \text{size}(B_v) \geqslant \max_{B_v \in \mathcal{B}_{\text{inner}}} \text{size}(B_v)$. We prove this by contradiction. Hence, we assume $\min_{B_v \in \mathcal{B}} \text{size}(B_v) < \max_{B_v \in \mathcal{B}_{\text{inner}}} \text{size}(B_v)$ and we define $B := \arg \min_{B_v \in \mathcal{B}} \text{size}(B_v)$ and $B' := \arg \max_{B_v \in \mathcal{B}_{\text{inner}}} \text{size}(B_v)$. Note that for any cube $B_v \in \mathcal{B}$ either $B_v$ itself is in $\mathcal{B}_{\text{inner}} \cup \mathcal{B}_{\text{leaf}}$ or $B_v = \text{parent}(B_w)$ for some cube $B_w \in \mathcal{B}_{\text{inner}} \cup \mathcal{B}_{\text{leaf}}$. We now make the following case distinction.

**Case I:** $B = \text{parent}(B_w)$ for some cube $B_w \in \mathcal{B}_{\text{inner}} \cup \mathcal{B}_{\text{leaf}}$. But this is an immediate contradiction since Algorithm 1 would have to split $B'$ before splitting $B$.

**Case II:** $B \in \mathcal{B}_{\text{inner}} \cup \mathcal{B}_{\text{leaf}}$. Because $B$ itself was put into $\mathcal{B}$ and not parent($B$), there exists a cube $B_w \in \mathcal{B}_{\text{inner}} \cup \mathcal{B}_{\text{leaf}}$ such that parent($B$) $\supsetneq$ parent($B_w$), which means $\text{size}(\text{parent}(B_w)) < \text{size}(\text{parent}(B))$. In order to complete the proof, it suffices to show that $\text{size}(\text{parent}(B_w)) \leqslant \text{size}(B)$. Indeed, since $B'$ has not been split by Algorithm 1 (because $B' \in \mathcal{B}_{\text{inner}}$) we know that $\text{size}(B') \leqslant \text{size}(\text{parent}(B_w))$. This inequality along with the inequality $\text{size}(\text{parent}(B_w)) \leqslant \text{size}(B)$ imply that $\text{size}(B') \leqslant \text{size}(B)$ which is in contradiction with $\text{size}(B) < \text{size}(B')$. To show that $\text{size}(\text{parent}(B_w)) \leqslant \text{size}(B)$ we consider the following two subcases. (i) parent($B$) is a degree-1 node. This means that parent($B$) corresponds to a cube that was split into a donut and the cube corresponding to $B$. Since the cube corresponding to $B_w$ must be completely inside the cube corresponding to parent($B$) (because $\text{size}(\text{parent}(B_w)) < \text{size}(\text{parent}(B))$) and a donut is empty we conclude that the cube corresponding to $B_w$ must be completely inside the cube corresponding to $B$. Hence, $\text{size}(\text{parent}(B_w)) \leqslant \text{size}(B)$. (ii) parent($B$) is not a degree-1 node. The inequality $\text{size}(\text{parent}(B_w)) < \text{size}(\text{parent}(B))$ along with the fact that parent($B$) is not a degree-1 node imply that $\text{size}(\text{parent}(B_w)) \leqslant \text{size}(B)$.

This completes the proof that LB is a correct lower bound. Next we prove that $R_Q$ is a weak $r$-packing for $r = \varepsilon \cdot \text{LB}/f(k)$. Observe that after the loop in lines 14–16, the set $\mathcal{B}_{\text{leaf}}$ is still a cube cover of $S_Q$. Moreover, each cube $B_v \in \mathcal{B}_{\text{leaf}}$ either contains a single point from $S_Q$ or its size is at most $r/\sqrt{d}$. Lemma 3 then implies that $R_Q$ is a weak $r$-packing for the desired value of $r$.

It remains to bound the size of $R_Q$. To this end we note that at each iteration of the loop in lines 3–9 the size of $\mathcal{B}_{\text{inner}} \cup \mathcal{B}_{\text{leaf}}$ increases by at most $2^d - 1$, so after the loop we have $|\mathcal{B}_{\text{inner}} \cup \mathcal{B}_{\text{leaf}}| \leqslant k2^{2d} + 2^d - 1$. The loop in lines 14–16 replaces each cube $B_v \in \mathcal{B}_{\text{inner}}$ by a number of smaller cubes. Since $\text{LB} = c \cdot \max_{B_v \in \mathcal{B}_{\text{inner}}} \text{size}(B_v)$ and $r = \varepsilon \cdot \text{LB}/f(k)$, each cube $B_v$ is replaced by only $O((f(k)2^d\sqrt{d}/(c\,\varepsilon))^d)$ smaller cubes. Since $d$ is a fixed constant, the total number of cubes we end up with (which is the same as the size of the $r$-packing) is $O(k(f(k)/(c\,\varepsilon))^d)$. ◀

Lemma 5, together with Lemma 2, establishes the correctness of our approach. To achieve a good running time, we need a few supporting data structures.

- We need a data structure that can answer the following queries: given a query box $Z$, find the deepest node $u$ in $\mathcal{T}(S)$ such that $Z \subseteq B_u$. With a *centroid-decomposition tree* $\mathcal{T}_{\mathrm{cd}}$ we can answer such queries in $O(\log n)$ time; see the full version for details.
- We need a data structure $\mathcal{D}$ that can answer the following queries on $S$: given a query box $Z$ and an integer $1 \leqslant i \leqslant d$, report a point in $S \cap Z$ with maximum $x_i$-coordinate, and one with minimum $x_i$-coordinate. It is possible to answer such queries in $O(\log^{d-1} n)$ time with a range tree (with fractional cascading), which uses $O(n \log^{d-1} n)$ storage. Note that this also allows us to compute the bounding box of $S \cap Z$ in $O(\log^{d-1} n)$ time. (In fact slightly better bounds are possible [19], but for simplicity we stick to using standard data structures.)

▶ **Lemma 6.** *Algorithm 1 runs in $O(k\,(f(k)/(c\,\varepsilon))^d + k\,((f(k)/(c\,\varepsilon))\log n)^{d-1})$ time.*

**Proof.** The two key observations in the proof are the following. First, when $B_v \nsubseteq Q$ we replace $v$ by the deepest node $u$ such that $\mathrm{bb}(S_Q \cap B_v) \subseteq B_u$, which implies at least two of the children of $u$ must contain a point in $S_Q$. From this we conclude that the number of iterations in Phase 1 is bounded by $k2^{2d}$. Second, we use the fact that in Phase 2 the computation of $\mathrm{bb}(S_Q \cap B_v)$ is only needed when $B_v \nsubseteq Q$. The complete proof is in the full version. ◀

This leads to the following theorem.

▶ **Theorem 7.** *Let $S$ be a set of $n$ points in $\mathbb{R}^d$ and let $\Phi$ be a $(c, f(k))$-regular cost function. Suppose we have an algorithm that solves the given clustering problem on a set of $m$ points in $T_{\mathrm{ss}}(m,k)$ time. Then there is a data structure that uses $O(n \log^{d-1} n)$ storage such that, for a query range $Q$ and query values $k \geqslant 2$ and $\varepsilon > 0$, we can compute a $(1+\varepsilon)$-approximate answer to a range-clustering query in time*

$$O\left(k\left(\frac{f(k)}{c\,\varepsilon}\cdot\log n\right)^{d-1} + T_{\mathrm{ss}}\left(k\left(\frac{f(k)}{c\,\varepsilon}\right)^d, k\right) + T_{\mathrm{expand}}(n,k)\right).$$

As an example application we consider $k$-center queries in the plane. (The result for rectilinear 2-center queries is actually inferior to the exact solution presented later.)

▶ **Corollary 8.** *Let $S$ be a set of $n$ points in $\mathbb{R}^2$. There is a data structure that uses $O(n \log n)$ storage such that, for a query range $Q$ and query values $k \geqslant 2$ and $\varepsilon > 0$, we can compute a $(1+\varepsilon)$-approximate answer to a $k$-center query within the following bounds:*
  - **(i)** *for the rectilinear case with $k = 2, 3$, the query time is $O((1/\varepsilon)\log n + 1/\varepsilon^2)$;*
  - **(ii)** *for the rectilinear case with $k = 4, 5$, the query time is $O((1/\varepsilon)\log n + (1/\varepsilon^2)\mathrm{polylog}(1/\varepsilon))$;*
  - **(iii)** *for the Euclidean case with $k = 2$, the expected query time is $O((1/\varepsilon)\log n + (1/\varepsilon^2)\cdot \log^2(1/\varepsilon))$;*
  - **(iv)** *for the rectilinear case with $k > 5$ and the Euclidean case with $k > 2$ the query time is $O((k/\varepsilon)\log n + (k/\varepsilon)^{O(\sqrt{k})})$.*

**Proof.** Recall that the cost function for the $k$-center problem is $(1/(2\sqrt{d}), 1)$-regular for the rectilinear case and $(1/2, 1)$-regular for the Euclidean case. We now obtain our results by plugging in the appropriate algorithms for the single-shot version. For (i) we use the linear-time algorithm of Hoffmann [16], for (ii) we use the $O(n \cdot \mathrm{polylog} n)$ algorithm of Sharir and Welzl [23], for (iii) we use the $O(n \log^2 n)$ randomized algorithm of Eppstein [12], for (iv) we use the $n^{O(\sqrt{k})}$ algorithm of Agarwal and Procopiuc [3]. ◀

## 3    Approximate Capacitated $k$-Center Queries

In this section we study the capacitated variant of the rectilinear $k$-center problem in the plane. In this variant we want to cover a set $S$ of $n$ points in $\mathbb{R}^2$ with $k$ congruent squares of minimum size, under the condition that no square is assigned more than $\alpha \cdot n/k$ points, where $\alpha > 1$ is a given constant. For a capacitated rectilinear $k$-center query this means we want to assign no more than $\alpha \cdot |S_Q|/k$ points to each square. Our data structure will report a $(1 + \varepsilon, 1 + \delta)$-approximate answer to capacitated rectilinear $k$-center queries: given a query range $Q$, a natural number $k \geqslant 2$, a constant $\alpha > 1$, and real numbers $\varepsilon, \delta > 0$, it computes a set $\mathcal{C} = \{b_1, \ldots, b_k\}$ of congruent squares such that:

- Each $b_i$ can be associated to a subset $C_i \subseteq S_Q \cap b_i$ such that $\{C_1, \ldots, C_k\}$ is a $k$-clustering of $S_Q$ and $|C_i| \leqslant (1 + \delta)\alpha \cdot |S_Q|/k$; and
- The size of the squares in $\mathcal{C}$ is at most $(1+\varepsilon)\cdot\mathrm{OPT}_k(S_Q, \alpha)$, where $\mathrm{OPT}_k(S_Q, \alpha)$ is the value of an optimal solution to the problem on $S_Q$ with capacity upper bound $U_Q := \alpha \cdot |S_Q|/k$.

Thus we allow ourselves to violate the capacity constraint by a factor $1 + \delta$.

To handle the capacity constraints, it is not sufficient to work with $r$-packings – we also need $\delta$-approximations. Let $P$ be a set of points in $\mathbb{R}^2$. A *$\delta$-approximation* of $P$ with respect to rectangles is a subset $A \subseteq P$ such that for any rectangle $\sigma$ we have

$$\big| \, |P \cap \sigma|/|P| \; - \; |A \cap \sigma|/|A| \big| \leqslant \delta.$$

From now on, whenever we speak of $\delta$-approximations, we mean $\delta$-approximations with respect to rectangles. Our method will use a special variant of the capacitated $k$-center problem, where we also have points that must be covered but do not count for the capacity:

▶ **Definition 9.** Let $R \cup A$ be a point set in $\mathbb{R}^2$, $k \geqslant 2$ a natural number, and $U$ a capacity bound. The *0/1-weighted capacitated k-center problem* in $\mathbb{R}^2$ is to compute a set $\mathcal{C} = \{b_1, \ldots, b_k\}$ of congruent squares of minimum size where each $b_i$ is associated to a subset $C_i \subseteq (R \cup A) \cap b_i$ such that $\{C_1, \ldots, C_k\}$ is a $k$-clustering of $R \cup A$ and $|C_i \cap A| \leqslant U$.

For a square $b$, let $\mathrm{expand}(b, r)$ denote the square $b$ expanded by $r$ on each side (so its radius in the $L_\infty$-metric increases by $r$). Let 0/1-WEIGHTEDKCENTER be an algorithm for the single-shot capacitated rectilinear $k$-center problem. Our query algorithm is as follows.

$\textsc{CapacitatedKCenterQuery}(k, Q, \alpha, \varepsilon, \delta)$

1: Compute a lower bound LB on $\mathrm{OPT}_k(S_Q)$.
2: Set $r := \varepsilon \cdot \mathrm{LB}/f(k)$ and compute a weak $r$-packing $R_Q$ on $S_Q$.
3: Set $\delta_Q := \delta/16k^3$ and compute a $\delta_Q$-approximation $A_Q$ on $S_Q$.
4: Set $U := (1+\delta/2)\cdot\alpha\cdot|A_Q|/k$ and $\mathcal{C} := 0/1\text{-}\textsc{WeightedKCenter}(R_Q\cup A_Q, k, U)$.
5: $\mathcal{C}^* := \{\mathrm{expand}(b, r) : b \in \mathcal{C}\}$.
6: **return** $\mathcal{C}^*$.

Note that the lower bound computed in Step 1 is a lower bound on the uncapacitated problem (which is also a lower bound for the capacitated problem). Hence, for Steps 1 and 2 we can use the algorithm from the previous section. How Step 3 is done will be explained later. First we show that the algorithm gives a $(1 + \varepsilon, 1 + \delta)$-approximate solution. We start by showing that we get a valid solution that violates the capacity constraint by at most a factor $(1 + \delta)$. (See the full version for a proof.)

▶ **Lemma 10.** *Let $\{b_1, \ldots, b_k\} := \mathcal{C}^*$ be the set of squares computed in Step 5. There exists a partition $\{C_1, \ldots, C_k\}$ of $S_Q$ such that $C_i \subseteq b_i$ and $|C_i| \leqslant (1 + \delta) \cdot U_Q$ for each $1 \leqslant i \leqslant k$, and such a partition can be computed in $O(k^2 + n \log n)$ time.*

We also need to prove that we get a $(1 + \varepsilon)$-approximate solution. For this it suffices to show that an optimal solution $\mathcal{C}_{\mathrm{opt}}$ to the problem on $S_Q$ is a valid solution on $R_Q \cup A_Q$. We can prove this by a similar approach as in the proof of the previous lemma.

▶ **Lemma 11.** *The size of the squares in $\mathcal{C}^*$ is at most $(1 + \varepsilon) \cdot \mathrm{OPT}_k(S_Q, \alpha)$.*

To make CAPACITATEDKCENTERQUERY run efficiently, we need some more supporting data structures. In particular, we need to quickly compute a $\delta_Q$-approximation within our range $Q$. We use the following data structures.

- We compute a collection $A_1, \ldots, A_{\log n}$ where $A_i$ is a $(1/2^i)$-approximation on $S$ using the algorithm of Phillips [21]. This algorithm computes, given a planar point set $P$ of size $m$ and a parameter $\delta$, a $\delta$-approximation of size $O((1/\delta) \log^4(1/\delta) \cdot \mathrm{polylog}(\log(1/\delta)))$ in time $O((n/\delta^3) \cdot \mathrm{polylog}(1/\delta))$. We store each $A_i$ in a data structure for orthogonal range-reporting queries. If we use a range tree with fractional cascading, the data structure uses $O(|A_i| \log |A_i|)$ storage and we can compute all the points in $A_i \cap Q$ in time $O(\log n + |A_i \cap Q|)$.
- We store $S$ in a data structure for orthogonal range-counting queries. There is such a data structure that uses $O(n)$ storage and such that queries take $O(\log n)$ time [9].

We can now compute a $\delta_Q$-approximation for $S_Q$ as follows.

> DELTAAPPROX$(Q, \delta_Q)$
> 1: Find the smallest value for $i$ such that $\frac{1}{2^i} \leqslant \frac{\delta_Q}{4} \frac{|S_Q|}{|S|}$, and compute $A := Q \cap A_i$.
> 2: Compute a $(\delta_Q/2)$-approximation $A_Q$ on $A$ using the algorithm by Phillips [21].
> 3: **return** $A_Q$.

▶ **Lemma 12.** DELTAAPPROX$(Q, \delta_Q)$ *computes a $\delta_Q$-approximation of size $O\left(\frac{1}{\delta_Q} \mathrm{polylog} \frac{1}{\delta_Q}\right)$ on $S_Q$ in time $O\big((\log n/\delta_Q)^4 \mathrm{polylog}(\log n/\delta_Q)\big)$.*

The only thing left is now an algorithm $0/1$-WEIGHTEDKCENTER$(R_Q \cup A_Q, k, U)$ that solves the $0/1$-weighted version of the capacitated rectilinear $k$-center problem. Here we use the following straightforward approach. Let $m := |R_Q \cup A_Q|$. First we observe that at least one square in an optimal solution has points on opposite edges. Hence, to find the optimal size we can do a binary search over $O(m^2)$ values, namely the horizontal and vertical distances between any pair of points. Moreover, given a target size $s$ we can push all squares such that each has a point on its bottom edge and a point on its left edge. Hence, to test if there is a solution of a given target size $s$, we only have to test $O(m^{2k})$ sets of $k$ squares. To test such a set $\mathcal{C} = \{b_1, \ldots, b_k\}$ of squares, we need to check if the squares cover all points in $R_Q \cup A_Q$ and if we can assign the points to squares such that the capacity constraint is met. For the latter we need to solve a flow problem, which can be done in $O(m^2 k)$ time; see the full version. Thus each step in the binary search takes $O(m^{2k+2} k)$, leading to an overall time complexity for $0/1$-WEIGHTEDKCENTER$(R_Q \cup A_Q, k, U)$ of $O(m^{2k+2} k \log m)$, where $m = |R_Q \cup A_Q| = O\big(k/\varepsilon^2 + (1/\delta_Q) \mathrm{polylog}(1/\delta_Q)\big)$, where $\delta_Q = \Theta(\delta/k^3)$.

The following theorem summarizes the results in this section.

▶ **Theorem 13.** *Let $S$ be a set of $n$ points in $\mathbb{R}^2$. There is a data structure that uses $O(n \log n)$ storage such that, for a query range $Q$ and query values $k \geqslant 2$, $\varepsilon > 0$ and $\delta > 0$, we can compute a $(1 + \varepsilon, 1 + \delta)$-approximate answer to a rectilinear $k$-center query in $O^*((k/\varepsilon) \log n + ((k^3/\delta) \log n)^4 + (k/\varepsilon^2 + (k^3/\delta))^{2k+2})$ time, where the $O^*$-notation hides $O(\mathrm{polylog}(k/\delta))$ factors.*

Note that for constant $k$ and $\varepsilon = \delta$ the query time simplifies to $O^*((1/\varepsilon^4)\log^4 n + (1/\varepsilon)^{4k+4})$. Also note that the time bound stated in the theorem only includes the time to compute the set of squares defining the clustering. If we want to also report an appropriate assignment of points to the squares, we have to add an $O(k^2 + |S_Q|\log|S_Q|)$ term; see Lemma 10.

▶ Remark. The algorithm can be generalized to the rectilinear $k$-center problem in higher dimensions, and to the Euclidean $k$-center problem; we only need to plug in an appropriate $\delta$-approximation algorithm and an appropriate algorithm for the 0/1-weighted version of the problem.

## 4     Exact $k$-Center Queries in $\mathbb{R}^1$

In this section we consider $k$-center queries in $\mathbb{R}^1$. Here we are given a set $S$ of $n$ points in $\mathbb{R}^1$ that we wish to preprocess into a data structure such that, given a query interval $Q$ and a natural number $k \geqslant 2$, we can compute a set $\mathcal{C}$ of at most $k$ intervals of the same length that together cover all points in $S_Q := S \cap Q$ and whose length is minimum. We obtain the following result (complete proof in the full version).

▶ **Theorem 14.** *Let $S$ be a set of $n$ points in $\mathbb{R}^1$. There is a data structure that uses $O(n)$ storage such that, for a query range $Q$ and a query value $k \geqslant 2$, we can answer a rectilinear $k$-center query in $O(\min\{k^2\log^2 n, 3^k\log n\})$ time.*

**Proof Sketch.** We present two approaches, one with $O(k^2\log^2 n)$ query time and one with $O(3^k\log n)$ query time. Let $p_1, \ldots, p_n$ be the points in $S$, sorted from left to right.

The first approach uses a subroutine DECIDER which, given an interval $Q'$, a length $L$ and an integer $\ell \leqslant k$, can decide in $O(\ell \log n)$ time if all points in $S \cap Q'$ can be covered by $\ell$ intervals of length $L$. The global query algorithm then performs a binary search, using DECIDER as subroutine, to find a pair of points $p_i, p_{i+1} \in S_Q$ such that the first interval in an optimal solution covers $p_i$ but not $p_{i+1}$. Then an optimal solution is found recursively for $k-1$ clusters within the query interval $Q \cap [p_{i+1}, \infty)$.

The second approach searches for a point $q \in Q$ and value $\ell$ such that there is an optimal solution where $S_Q \cap (-\infty, q]$ is covered by $\ell$ intervals and $S_Q \cap [q, \infty)$ is covered by $k-\ell$ intervals. The efficiency depends on the interesting fact that there must be a *fair split point*, that is, a pair $(q, \ell)$ such that $\ell/k$ (the fraction of intervals used to the left of $q$) is proportional to the fraction of the length of $Q$ to the left of $q$.    ◀

## 5     Exact Rectilinear 2- and 3-Center Queries in $\mathbb{R}^2$

Suppose we are given a set $S = \{p_1, p_2, \ldots, p_n\}$ of $n$ points in $\mathbb{R}^2$ and an integer $k$. In this section we build a data structure $\mathcal{D}$ that stores the set $S$ and, given an orthogonal query rectangle $Q$, can be used to quickly find an optimal solution for the $k$-center problem on $S_Q$ for $k = 2$ or 3, where $S_Q := S \cap Q$.

**2-center queries**

Our approach for the case of $k = 2$ is as follows. We start by shrinking the query range $Q$ such that each edge of $Q$ touches at least one point of $S$. (The time for this step is subsumed by the time for the rest of the procedure.) It is well known that if we want to cover $S_Q$ by two squares $\sigma, \sigma'$ of minimum size, then $\sigma$ and $\sigma'$ both share a corner with $Q$ and these corners are opposite corners of $Q$. We say that $\sigma$ and $\sigma'$ are *anchored* at the corner they share with $Q$. Thus we need to find optimal solutions for two cases – $\sigma$ and $\sigma'$ are anchored at the
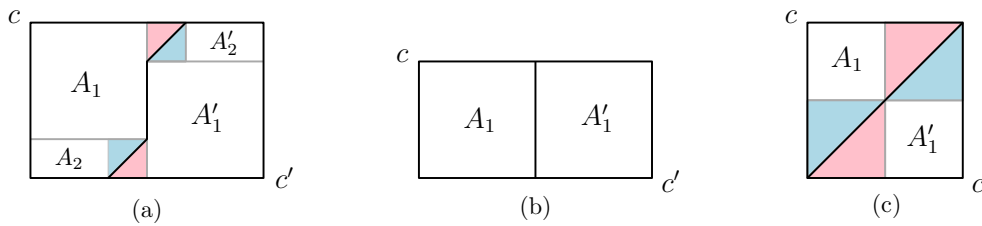
**Figure 1** Various types of $L_\infty$-bisectors. The bisectors and the boundaries of query regions are shown in black. (a): $Q$ is "fat". The regions $A_j, A'_j$ for $j = 1, 2$ are shown with text. (b): $Q$ is "thin". The regions $A_j$ and $A'_j$ for $j = 2, 3, 4$ are empty. (c): $Q$ is a square. The regions $A_j$ and $A'_j$ for $j = 2$ are empty. In both (a) and (c) regions $A_3, A'_3$ are colored in blue and $A_4, A'_4$ are colored in pink.

topleft and bottomright corner of $Q$, or at the topright and bottomleft corner – and return the better one. Let $c$ and $c'$ be the topleft and the bottomright corners of $Q$, respectively. In the following we describe how to compute two squares $\sigma$ and $\sigma'$ of minimum size that are anchored at $c$ and $c'$, respectively, and whose union covers $S_Q$. The topright/bottomleft case can then be handled in the same way.

First we determine the $L_\infty$-bisector of $c$ and $c'$ inside $Q$; see Figure 1. The bisector partitions $Q$ into regions $A$ and $A'$, that respectively have $c$ and $c'$ on their boundary. Obviously in an optimal solution (of the type we are focusing on), the square $\sigma$ must cover $S_Q \cap A$ and the square $\sigma'$ must cover $S_Q \cap A'$. To compute $\sigma$ and $\sigma'$, we thus need to find the points $q \in A$ and $q' \in A'$ with maximum $L_\infty$-distance to the corners $c$ and $c'$ respectively. To this end, we partition $A$ and $A'$ into subregions such that in each of the subregions the point with maximum $L_\infty$-distance to its corresponding corner can be found quickly via appropriate data structures discussed below. We assume w.l.o.g. that the $x$-span of $Q$ is at least its $y$-span. We begin by presenting the details of such a partitioning for Case (a) of Figure 1 – Cases (b) and (c) can be seen as special cases of Case (a).

As Figure 1 suggests, we partition $A$ and $A'$ into subregions. We denote these subregions by $A_j$ and $A'_j$, for $1 \leqslant j \leqslant 4$. From now on we focus on reporting the point $q \in S$ in $A$ with maximum $L_\infty$-distance to $c$; finding the furthest point from $c'$ inside $A'$ can be done similarly. Define four points $p(A_j) \in S$ for $1 \leqslant j \leqslant 4$ as follows.

- The point $p(A_1)$ is the point of $S_Q$ with maximum $L_\infty$-distance to $c$ in $A_1$. Note that this is either the point with maximum $x$-coordinate in $A_1$ or the point with minimum $y$-coordinate.
- The point $A_2$ is a bottommost point in $A_2$.
- The point $A_3$ is a bottommost point in $A_3$.
- The point $A_4$ is a rightmost point in $A_4$.

Clearly

$$q = \arg \max_{1 \leqslant j \leqslant 4} \{d_\infty(p(A_j), c)\}, \tag{1}$$

where $d_\infty(.)$ denotes the $L_\infty$-distance function.

**Data structure.** Our data structure now consists of the following components.

- We store $S$ in a data structure $\mathcal{D}_1$ that allows us to report the extreme points in the $x$-direction and in the $y$-direction inside a rectangular query range. For this we use the structure by Chazelle [9], which uses $O(n \log^\varepsilon n)$ storage and has $O(\log n)$ query time.
- We store $S$ in a data structure $\mathcal{D}_2$ with two components. The first component should answer the following queries: given a 45° query cone whose top bounding line is horizontal

and that is directed to the left – we obtain such a cone when we extend the region $A_4$ into an infinite cone –, report the rightmost point inside the cone. The second component should answer similar queries for cones that are the extension of $A_3$.

In the full version we describe a linear-size data structure for such queries that has $O(\log n)$ query time.

**Query procedure.** Given an axis-aligned query rectangle $Q$, we first (as already mentioned) shrink the query range so that each edge of $Q$ contains at least on point of $S$. Then compute the $L_\infty$-bisector of $Q$. Query $\mathcal{D}_1$ with $A_1$ and $A_2$, respectively, to get the points $p(A_1)$ and $p(A_2)$. Then query $\mathcal{D}_2$ with $u$ and $u'$ to get the points $p(A_3)$ and $p(A_4)$, where $u$ and $u'$ are respectively the bottom and the top intersection points of $L_\infty$-bisector of $Q$ and the boundary of $Q$. Among the at most four reported points, take the one with maximum $L_\infty$-distance the corner $c$. This is the point $q \in S_Q \cap A$ furthest from $c$.

Compute the point $q' \in S_Q \cap A'$ furthest from $c'$ in a similar fashion. Finally, report two minimum-size congruent squares $\sigma$ and $\sigma'$ anchored at $c$ and $c'$ and containing $q$ and $q'$, respectively.

Putting everything together, we end up with the following theorem.

▶ **Theorem 15.** *Let $S$ be a set of $n$ points in $\mathbb{R}^2$. For any fixed $\varepsilon > 0$, there is a data structure using $O(n \log^\varepsilon n)$ storage that can answer rectilinear 2-center queries in $O(\log n)$ time.*

▶ Remark. We note that the query time in Theorem 15 can be improved in the word-RAM model to $O(\log \log n)$ by using the range successor data structure of Zhou [24], and the point location data structure for orthogonal subdivisions by de Berg *et al.* [11].

### 3-center queries

Given a (shrunk) query range $Q$, we need to compute a set $\{\sigma_1, \sigma_2, \sigma_3\}$ of (at most) three congruent squares of minimal size whose union covers $S_Q$. It is easy to verify (and is well-known) that at least one of the squares in an optimal solution must be anchored at one of the corners of $Q$. Hence and w.l.o.g. we assume that $\sigma_1$ is anchored at one of the corners of $Q$. We try placing $\sigma_1$ in each corner of $Q$ and select the placement resulting in the best overall solution. Next we briefly explain how to find the best solution subject to placing $\sigma_1$ in the leftbottom corner of $Q$. The other cases are symmetric. We perform two separate binary searches; one will test placements of $\sigma_1$ such that its right side has the same $x$-coordinate as a point in $S$, the other will be on possible $y$-coordinates for the top side. During each of the binary searches, we compute the smallest axis-parallel rectangle $Q' \subseteq Q$ containing the points of $Q \setminus \sigma_1$ (by covering $Q \setminus \sigma_1$ with axis-parallel rectangles and querying for extreme points in these rectangles). We then run the algorithm for $k = 2$ on $Q'$. We need to ensure that this query ignores the points already covered by $\sigma_1$. For this, recall that for $k = 2$ we covered the regions $A$ and $A'$ by suitable rectangular and triangular ranges. We can now do the same, but we cover $A \setminus \sigma_1$ and $A' \setminus \sigma_1$ instead.

After the query on $Q'$, we compare the size of the resulting squares with the size of $\sigma_1$ to guide the binary search. The process stops as soon as the three sizes are the same or no further progress in the binary search can be made. Putting everything together, we end up with the following theorem.

▶ **Theorem 16.** *Let $S$ be a set of $n$ points in $\mathbb{R}^2$. For any fixed $\varepsilon > 0$, there is a data structure using $O(n \log^\varepsilon n)$ storage that can answer rectilinear 3-center queries in $O(\log^2 n)$ time.*

▶ Remark. Similar to Theorem 15, the query time in Theorem 16 can be improved in the word-RAM model of computation to $O(\log n \log \log n)$ time.

―――― **References** ――――

**1**  M. Abam, P. Carmi, M. Farshi, and M. Smid. On the power of the semi-separated pair decomposition. *Compututational Geometry: Theory and Applications*, 46:631–639, 2013.

**2**  P. K. Agarwal, R. Ben Avraham, and M. Sharir. The 2-center problem in three dimensions. *Compututational Geometry: Theory and Applications*, 46:734–746, 2013.

**3**  P. K. Agarwal and Cecilia M. Procopiuc. Exact and approximation algorithms for clustering. *Algorithmica*, 33:201–226, 2002.

**4**  Sunil Arya, David M. Mount, and Eunhui Park. Approximate geometric MST range queries. In *Proc. 36th International Symposium on Computational Geometry (SoCG)*, pages 781–795, 2015.

**5**  Peter Brass, Christian Knauer, Chan-Su Shin, Michiel H. M. Smid, and Ivo Vigan. Range-aggregate queries for geometric extent problems. In *Computing: The Australasian Theory Symposium 2013, CATS'13*, pages 3–10, 2013.

**6**  V. Capoyleas, G. Rote, and G. Woeginger. Geometric clusterings. *Journal of Algorithms*, 12:341–356, 1991.

**7**  T. M. Chan. Geometric applications of a randomized optimization technique. *Discrete & Compututational Geometry*, 22:547–567, 1999.

**8**  T. M. Chan. More planar two-center algorithms. *Compututational Geometry: Theory and Applications*, 13:189–198, 1999.

**9**  Bernard Chazelle. A functional approach to data structures and its use in multidimensional searching. *SIAM Journal on Computing*, 17:427–462, 1988.

**10**  A. W. Das, P. Gupta, K. Kothapalli, and K. Srinathan. On reporting the $L_1$-metric closest pair in a query rectangle. *Information Processing Letters*, 114:256–263, 2014.

**11**  Mark de Berg, Marc van Kreveld, and Jack Snoeyink. Two- and three-dimensional point location in rectangular subdivisions. *Journal of Algorithms*, 18:256–277, 1995.

**12**  D. Eppstein. Faster construction of planar two-centers. In *Proc. 8th Annual ACM-SIAM Symposiun on Discrete Algorithms (SODA)*, pages 131–138, 1997.

**13**  P. Gupta, R. Janardan, Y. Kumar, and M. Smid. Data structures for range-aggregate extent queries. *Compututational Geometry: Theory and Applications*, 47:329–347, 2014.

**14**  Sariel Har-Peled. *Geometric Approximation Algorithms*, volume 173 of *Mathematical surveys and monographs*. American Mathematical Society, 2011.

**15**  Sariel Har-Peled and Soham Mazumdar. On coresets for $k$-means and $k$-median clustering. In *Proc. 36th Annual ACM Symposium on Theory of Computing (STOC)*, pages 291–300, 2004.

**16**  M. Hoffmann. A simple linear algorithm for computing rectilinear 3-centers. *Compututational Geometry: Theory and Applications*, 31:150–165, 2005.

**17**  R. Z. Hwang, R. Lee, and R. C. Chang. The generalized searching over separators strategy to solve some NP-hard problems in subexponential time. *Algorithmica*, 9:398–423, 1993.

**18**   S. Khare, J. Agarwal, N. Moidu, and K. Srinathan. Improved bounds for smallest enclosing disk range queries. In *Proc. 26th Canadian Conference on Computational Geometry (CCCG)*, 2014.

**19**   H.-P. Lenhof and M. H. M. Smid. Using persistent data structures for adding range restrictions to searching problems. *Theoretical Informatics and Applications*, 28:25–49, 1994.

**20**   Yakov Nekrich and Michiel H. M. Smid. Approximating range-aggregate queries using coresets. In *Proc. 22nd Canadian Conference on Computational Geometry (CCCG)*, pages 253–256, 2010.

**21**   Jeff M. Phillips. Algorithms for $\varepsilon$-approximations of terrains. In *Proc. 35th International Colloquium on Automata, Languages, and Programming (ICALP)*, pages 447–458, 2008.

**22**   M. Sharir. A near-linear time algorithm for the planar 2-center problem. *Discrete & Compututational Geometry*, 18:125–134, 1997.

**23**   M. Sharir and E. Welzl. Rectilinear and polygonal $p$-piercing and $p$-center problems. In *Proc. 12th International Symposium on Computational Geometry (SoCG)*, pages 122–132, 1996.

**24**   Gelin Zhou. Two-dimensional range successor in optimal time and almost linear space. *Information Processing Letters*, 116:171–174, 2016.