

## Approximating Minimization Diagrams and Generalized Proximity Search

Sariel Har-Peled\*  
 Dept. of Computer Science  
 University of Illinois,  
 Urbana, IL 61801, USA  
 Email: sariel@uiuc.edu

Nirman Kumar\*  
 Dept. of Computer Science  
 University of Illinois,  
 Urbana, IL 61801, USA  
 Email: nkumar5@illinois.edu

**Abstract**—We investigate the classes of functions whose minimization diagrams can be approximated efficiently in  $\mathbb{R}^d$ . We present a general framework and a data-structure that can be used to approximate the minimization diagram of such functions. The resulting data-structure has near linear size and can answer queries in logarithmic time. Applications include approximating the Voronoi diagram of (additively or multiplicatively) weighted points. Our technique also works for more general distance functions, such as metrics induced by convex bodies, and the nearest furthest-neighbor distance to a set of point sets. Interestingly, our framework works also for distance functions that do not obey the triangle inequality. For many of these functions no near-linear size approximation was known before.

## I. INTRODUCTION

Let  $\mathcal{F} = \{f_i : \mathbb{R}^d \rightarrow \mathbb{R} \mid i = 1, \dots, n\}$ , be a set of functions. The minimization diagram of  $\mathcal{F}$  is the function  $f_{\min}(\mathbf{q}) = \min_{i=1, \dots, n} f_i(\mathbf{q})$ , for any  $\mathbf{q} \in \mathbb{R}^d$ . By viewing the graphs of these functions as manifolds in  $\mathbb{R}^{d+1}$ , the graph of the minimization diagram, also known as the *lower envelope* of  $\mathcal{F}$ , is the manifold that can be viewed from an observer at  $-\infty$  on the  $x_{d+1}$  axis. Given a set of functions  $\mathcal{F}$  as above, many problems in Computational Geometry can be viewed as computing the minimization diagram; that is, one preprocesses  $\mathcal{F}$ , and given a query point  $\mathbf{q}$ , one needs to compute  $f_{\min}(\mathbf{q})$  quickly. This typically requires  $n^{O(d)}$  space if one is interested in logarithmic query time. If one is restricted to using linear space, then the query time deteriorates to  $O(n^{1-O(1/d)})$  [2], [3]. There is substantial work on bounding the complexity of the lower envelope in various cases, how to compute it efficiently, and performing range search on them; see the book by Sharir and Agarwal [4].

*Nearest neighbor.*

One natural problem that falls into this framework is the nearest neighbor (NN) search problem. Here, given a set  $P$  of  $n$  data points in a metric space  $\mathcal{X}$ , we need to preprocess  $P$ , such that given a query point  $\mathbf{q} \in \mathcal{X}$ , one can find (quickly) the point  $\mathbf{n}_q \in P$  closest to  $\mathbf{q}$ . Nearest neighbor search is a fundamental task used in numerous domains including machine learning, clustering, document retrieval, databases, statistics, and many others.

\*Work on this paper was partially supported by NSF AF award CCF-0915984, and NSF AF award CCF-1217462. The full version of this paper is available at [1].

To see the connection to lower envelopes, consider a set of data points  $P = \{p_1, \dots, p_n\}$  in  $\mathbb{R}^d$ . Next, consider the set of functions  $\mathcal{F} = \{f_1, \dots, f_n\}$ , where  $f_i(\mathbf{q}) = \|\mathbf{q} - p_i\|$ , for  $i = 1, \dots, n$ . The graph of  $f_i$  is the set of points  $\{(\mathbf{q}, f_i(\mathbf{q})) \mid \mathbf{q} \in \mathbb{R}^d\}$  (which is a cone in  $\mathbb{R}^{d+1}$  with apex at  $(p_i, 0)$ ). Clearly the NN problem is to evaluate the minimization diagram of the functions at a query point  $\mathbf{q}$ .

More generally, given a set of  $n$  functions, one can think of the minimization diagram defining a “distance function”, by analogy with the above. The distance of a query point here is simply the “height” of the lower envelope at that point.

*Exact nearest neighbor.*

The exact nearest neighbor problem has a naive linear time algorithm without any preprocessing. However, by doing some nontrivial preprocessing, one can achieve a sub-linear query time. In  $\mathbb{R}^d$ , this is facilitated by answering point location queries using a Voronoi diagram [5]. However, this approach is only suitable for low dimensions, as the complexity of the Voronoi diagram is  $\Theta(n^{\lceil d/2 \rceil})$  in the worst case. Specifically, Clarkson [6] showed a data-structure with query time  $O(\log n)$  time, and  $O(n^{\lceil d/2 \rceil + \delta})$  space, where  $\delta > 0$  is a prespecified constant (the  $O(\cdot)$  notation here hides constants that are exponential in the dimension). One can trade-off the space used and the query time [7]. Meiser [8] provided a data-structure with query time  $O(d^5 \log n)$  (which has polynomial dependency on the dimension), where the space used is  $O(n^{d+\delta})$ . These solutions are impractical even for data-sets of moderate size if the dimension is larger than two.

*Approximate nearest neighbor.*

In typical applications however, it is usually sufficient to return an *approximate nearest neighbor (ANN)*. Given an  $\varepsilon > 0$ , a  $(1 + \varepsilon)$ -ANN, to a query point  $\mathbf{q}$ , is a point  $y \in P$ , such that

$$\|\mathbf{q} - y\| \leq (1 + \varepsilon) \|\mathbf{q} - \mathbf{n}_q\|,$$

where  $\mathbf{n}_q \in P$  is the nearest neighbor to  $\mathbf{q}$  in  $P$ . Considerable amount of work was done on this problem, see [9] and references therein.

Indyk and Motwani showed that ANN queries can be reduced to a small number of near neighbor queries [10], [11]. For high dimensional Euclidean spaces, they used locality sensitive hashing to solve the near neighbor problem and provided a data-structure that answers ANN queries in time (roughly)  $\tilde{O}(n^{1/(1+\varepsilon)})$  with preprocessing time and space

$\tilde{O}(n^{1+1/(1+\varepsilon)})$ ; here the  $\tilde{O}(\cdot)$  hides terms polynomial in  $\log n$  and  $1/\varepsilon$ . This was improved to  $\tilde{O}(n^{1/(1+\varepsilon)^2})$  query time, and preprocessing time and space  $\tilde{O}(n^{1+1/(1+\varepsilon)^2})$  [12]. These bounds are near optimal [13].

In low dimensions (i.e.,  $\mathbb{R}^d$  for small  $d$ ), one can use linear space (independent of  $\varepsilon$ ) and get ANN query time  $O(\log n + 1/\varepsilon^{d-1})$  [14], [15]. The trade-off for this logarithmic query time is of course an exponential dependence on  $d$ . Interestingly, for this data-structure, the approximation parameter  $\varepsilon$  is not prespecified during the construction; one needs to provide it only during the query. An alternative approach is to use Approximate Voronoi Diagrams (AVD), introduced by Har-Peled [16], which is a partition of space into regions, of near-linear total complexity, typically with a representative point for each region that is an ANN for any point in the region. In particular, Har-Peled showed that there is such a decomposition of size  $O((n/\varepsilon^d) \log^2 n)$ , such that ANN queries can be answered in  $O(\log n)$  time. Arya and Malamatos [17] showed how to build AVD's of linear complexity (i.e.,  $O(n/\varepsilon^d)$ ). Their construction uses Well-Separated Pair Decomposition [18]. Further trade-offs between query time and space for AVD's were studied by Arya *et al.* [19].

*Generalized distance functions: motivation.*

The algorithms for approximate nearest neighbor, extend to various metrics in  $\mathbb{R}^d$ , for example the well known  $\ell_p$  metrics. In particular, previous constructions of AVD's extend to  $\ell_p$  metrics [16], [17] as well. However, these constructions fail even for a relatively simple and natural extension; specifically, multiplicative weighted Voronoi diagrams. Here, every site  $\mathbf{p}$ , in the given point set  $P$ , has a weight  $\omega_{\mathbf{p}}$ , and the “distance” of a query point  $\mathbf{q}$  to  $\mathbf{p}$  is  $f_{\mathbf{p}}(\mathbf{q}) = \omega_{\mathbf{p}} \|\mathbf{q} - \mathbf{p}\|$ . The function  $f_{\mathbf{p}}$  is the natural distance function induced by  $\mathbf{p}$ . As with ordinary Voronoi diagrams, one can define the weighted Voronoi diagram as a partition of space into disjoint regions, one for each site  $\mathbf{p}$ , such that in the region for  $\mathbf{p}$  the function  $f_{\mathbf{p}}$  is the one realizing the minimum among all the functions induced by the points of  $P$ . It is known that, even in the plane, multiplicative Voronoi diagrams can have quadratic complexity, and the minimizing distance function usually does not obey the triangle inequality. Intuitively, such multiplicative Voronoi diagrams can be used to model facilities where the price of delivery to a client depends on the facility and the distance. Of course, this is only one possible distance function, and there are many other such functions that are of interest (e.g., multiplicative, additive, etc.).

*When fast proximity and small space is not possible.*

Consider a set of segments in the plane, and we are interested in the nearest segment to a query point. Given  $n$  such segments and  $n$  such query points, this is an extension of Hopcroft's problem, which requires only to decide if there is any of the given points on any of the segments. There are lower bounds (in reasonable models) that show that Hopcroft's problem cannot be solved faster than  $\Omega(n^{4/3})$  time [20]. This implies that no multiplicative-error approximation for proximity search in this case is possible, if one

insists on near linear preprocessing, and logarithmic query time.

*When is fast ANN possible.*

So, consider a set of geometric objects where each one of them induces a natural distance function, measuring how far a point in space is from this object. Given such a collection of functions, the nearest neighbor for a query point is simply the function that defines the lower envelope “above” the query point (i.e., the object closest to the query point under its distance function). Clearly, this approach allows a generalization of the proximity search problem. In particular, the above question becomes, for what classes of functions, can the lower envelope be approximated up to  $(1 + \varepsilon)$ -multiplicative error, in logarithmic time? Here the preprocessing space used by the data structure should be near linear.

*I-A. Our results*

We characterize the conditions that are sufficient to efficiently approximate the minimization diagram of functions. Using this framework, one can quickly evaluate the lower envelope approximately for large classes of functions that arise naturally from proximity problems. Our data-structure can be constructed in near linear time, uses near linear space, and answers proximity queries in logarithmic time (for constant dimension  $d$ ). Our framework is quite general and should be applicable to many distance functions, and in particular we present the following specific cases where the new data-structure can be used:

- (A) **Multiplicative Voronoi diagrams.** Given a set of points  $P$ , where the  $i$ th point  $\mathbf{p}_i$  has associated weight  $w_i > 0$ , for  $i = 1, \dots, n$ , consider the functions  $f_i(\mathbf{q}) = w_i \|\mathbf{q} - \mathbf{p}_i\|$ . The minimization diagram for this set of functions, corresponds to the multiplicative weighted Voronoi diagram of the points. The approach of Arya and Malamatos [17] to construct AVD's using WSPD's fails for this problem, as that construction relies on the triangle inequality that the regular Euclidean distance possesses, which does not hold in this case. We provide a near linear space AVD construction for this case. We are unaware of any previous results on AVD for multiplicatively weighted Voronoi diagrams.
- (B) **Minkowski norms of fat convex bodies.** Given a bounded symmetric convex body  $C$  centered at the origin, it defines a natural metric; that is, for points  $\mathbf{u}$  and  $\mathbf{v}$  their distance, as induced by  $C$ , denoted by  $\|\mathbf{u} - \mathbf{v}\|_C$ , is the minimum  $x$  such that  $x\mathbf{C} + \mathbf{u}$  contains  $\mathbf{v}$ . So, given a set of  $n$  data points  $P = \{\mathbf{p}_1, \dots, \mathbf{p}_n\}$  and  $n$  centrally symmetric and bounded convex bodies  $C_1, \dots, C_n$ , we define  $f_i(\mathbf{q}) = \|\mathbf{p}_i - \mathbf{q}\|_{C_i}$ , for  $i = 1, \dots, n$ . Since each point induces a distance by a different convex body, this collection no longer defines a metric, and this makes the problem significantly more challenging. In particular, existing techniques for AVD and ANN cannot be readily applied. Intuitively, the fatness of the associated convex bodies turns out to be sufficient to approximate the associated distance function, see Section IV-B.
- (C) **Nearest furthest-neighbor.** Consider a situation where the given input is uncertain; specifically, for the

$i$ th point we are given a set of points  $P_i \subseteq \mathbb{R}^d$  where it might lie (the reader might consider the case where the  $i$ th point randomly chooses its location out of the points of  $P_i$ ). There is a growing interest in how to handle such inputs, as real world measurements are fraught with uncertainty, see [21], [22], [23], [24] and references therein. In particular, in the worst case, the distance of the query point  $\mathbf{q}$  to the  $i$ th point, is the distance from  $\mathbf{q}$  to the furthest-neighbor of  $\mathbf{q}$  in  $P_i$ ; that is,  $\mathcal{F}_i(\mathbf{q}) = \max_{\mathbf{p} \in P_i} \|\mathbf{q} - \mathbf{p}\|$ . Thus, in the worst case, the nearest point to the query is  $\mathcal{F}(\mathbf{q}) = \min_i \mathcal{F}_i(\mathbf{q})$ . Using our framework we can approximate this function efficiently, using space  $\tilde{O}(n)$ , and providing logarithmic query time. Note that surprisingly, the space requirement is independent of the original input size, and only depends on the number of uncertain points.

*Paper organization.*

In Section II we define our framework and prove some basic properties. Since we are trying to make our framework as inclusive as possible, its description is somewhat abstract. In Section III, we describe the construction of the AVD and its associated data-structure. We describe in Section IV some specific cases where the new AVD construction can be used. We conclude in Section V.

## II. PRELIMINARIES

For the sake of simplicity of exposition, throughout the paper we assume that all the “action” takes place in the unit cube  $[0, 1]^d$ . Among other things, this implies that all the queries are in this region. This can always be guaranteed by an appropriate scaling and translation of space. The scaling and translation, along with the conditions on functions in our framework, implies that outside the unit cube the approximation to the lower envelope can be obtained in constant time.

### II-A. Notations and basic definitions

Given  $\mathbf{q} \in \mathbb{R}^d$  and  $P \subseteq \mathbb{R}^d$  a non-empty closed set, the **distance** of  $\mathbf{q}$  to  $P$  is  $d(\mathbf{q}, P) = \min_{x \in P} \|\mathbf{q} - x\|$ . For a number  $\ell > 0$ , the **grid** of side-length  $\ell$ , denoted by  $G_\ell$ , is the natural tiling of  $\mathbb{R}^d$ , with cubes of side-length  $\ell$  (i.e., with a vertex at the origin). A cube  $\square$  is **canonical** if it belongs to  $G_\ell$ ,  $\ell$  is a power of 2, and  $\square \subseteq [0, 1]^d$ . Informally, a canonical cube (or cell) is a region that might correspond to a cell in a quadtree having the unit cube as the root region.

**Definition II.1.** To approximate a set  $X \subseteq [0, 1]^d$ , up to distance  $r$ , consider the set  $G_{\approx r}(X)$  of all the canonical grid cells of  $G_\ell$  that have a non-empty intersection with  $X$ , where  $\ell = 2^{\lfloor \log_2(r/\sqrt{d}) \rfloor}$ . Let  $\cup G_{\approx r}(X) = \cup_{\square \in G_{\approx r}(X)} \square$ , denote the union of cubes of  $G_{\approx r}(X)$ .

Observe that  $X \subseteq \cup G_{\approx r}(X) \subseteq X \oplus B(0, r)$ , where  $\oplus$  denotes the Minkowski sum, and  $B(0, r)$  is the ball of radius  $r$  centered at the origin.

**Definition II.2.** For  $\ell \geq 0$  and a function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$ , the  $\ell$  **sublevel set** of  $f$  is the set  $f_{\leq \ell} = \{\mathbf{p} \in \mathbb{R}^d \mid f(\mathbf{p}) \leq \ell\}$ . For a set of functions  $\mathcal{F}$ , let  $\mathcal{F}_{\leq \ell} = \cup_{f \in \mathcal{F}} f_{\leq \ell}$ .

**Definition II.3.** Given a function  $f$  and  $\mathbf{q} \in \mathbb{R}^d$  their distance is  $d(\mathbf{q}, f) = f(\mathbf{q})$ . Given two functions  $f$  and  $g$ , their **distance**  $d(f, g)$  is the minimum  $l \geq 0$  such that  $f_{\leq l} \cap g_{\leq l} \neq \emptyset$ . Similarly, for two sets of function,  $\mathcal{F}$  and  $\mathcal{G}$ , their **distance** is

$$d(\mathcal{F}, \mathcal{G}) = \min_{f \in \mathcal{F}, g \in \mathcal{G}} d(f, g).$$

The distance function behaves to some extent like one would expect an usual metric to behave: (i)  $d(f, g)$  always exists, and (ii) (symmetry)  $d(f, g) = d(g, f)$ . Also, we have  $f_{\leq d(f, g)} \neq \emptyset$ . We extend the above definition to sets of functions. Note that the triangle inequality does not hold for  $d(\cdot, \cdot)$ .

**Observation II.4.** Suppose that  $f$  and  $g$  are two functions such that  $d(f, g) > 0$  and  $\mathbf{q} \in \mathbb{R}^d$ . Then,  $\max(d(\mathbf{q}, f), d(\mathbf{q}, g)) \geq d(f, g)$ .

**Definition II.5.** Let  $B_1, B_2, \dots, B_m$  be  $n$  connected, nonempty sets in  $\mathbb{R}^d$ . This collection of sets is **connected** if  $\cup_i B_i$  is connected.

### II-A1. Sketches

A key idea underlying our approach is that is that any set of functions of interest should look like a single (or a small number of functions) from “far” enough. Indeed, given a set of points  $P \subseteq \mathbb{R}^d$ , they look like a single point (as far as distance), if the distance from  $\mathcal{CH}(P)$  is at least  $2\text{diam}(P)/\varepsilon$ .

**Definition II.6 (cl( $\mathcal{F}$ )).** Given a set of functions  $\mathcal{G}$ , if  $\mathcal{G}$  contains a single function then the **connectivity level**  $\text{cl}(\mathcal{G})$  is 0; otherwise, it is the minimum  $\ell \geq 0$ , such that the collection of sets  $f_{\leq \ell}$  for  $f \in \mathcal{G}$  is connected, see Definition II.5.

**Remark II.7.** It follows from Definition II.6 that at level  $\ell = \text{cl}(\mathcal{G})$ , each of the sets  $f_{\leq \ell}$  for  $f \in \mathcal{G}$  are nonempty and connected and further their union  $\mathcal{G}_{\leq \ell}$  is also connected. This can be relaxed to require that the intersection graph of the sets  $f_{\leq \ell}$  for  $f \in \mathcal{G}$  is connected (this also implies they are nonempty). Notice that, if at level  $\ell$ , the sublevel sets are connected, then the relaxed definition is equivalent to Definition II.6. However, the relaxed definition introduces more technical baggage, and for all the interesting applications we have, the sublevel sets  $f_{\leq y}$  are connected at all levels  $y$  at which they are nonempty. Therefore, in the interest of brevity, and to keep the presentation simple, we mandate that the sublevel sets be connected at  $\ell$ . In fact, it would not harm to assume that the sublevel sets are connected whenever nonempty.

**Definition II.8.** Given a set of functions  $\mathcal{G}$  and  $\delta \geq 0, y_0 \geq 0$ , a  $(\delta, y_0)$ -**sketch** for  $\mathcal{G}$  is a (hopefully small) subset  $\mathcal{H} \subseteq \mathcal{G}$ , such that  $\mathcal{G}_{\leq y} \subseteq \mathcal{H}_{\leq (1+\delta)y}$ , for all  $y \geq y_0$ .

It is easy to see that for any  $\mathcal{G}, \delta \geq 0, y_0 \geq 0$ , if  $\mathcal{H} \subseteq \mathcal{G}$  is a  $(\delta, y_0)$ -sketch, then for any  $\delta' \geq \delta, y_0' \geq y_0, \mathcal{H}' \supseteq \mathcal{H}$  it is true that  $\mathcal{H}'$  is a  $(\delta', y_0')$ -sketch for  $\mathcal{G}$ . Trivially, for any  $\delta \geq 0, y_0 \geq 0$ , it is true that  $\mathcal{H} = \mathcal{G}$  is a  $(\delta, y_0)$ -sketch.

### II-B. Conditions on the functions

We require that the set of functions under consideration satisfy the following conditions.

- (P1) **Compactness.** For any  $y \geq 0$  and  $i = 1, \dots, n$ , the set  $(f_i)_{\leq y}$  is compact.
- (P2) **Bounded growth.** For any  $f \in \mathcal{F}$ , there is a function  $\lambda_f : \mathbb{R}^+ \rightarrow \mathbb{R}^+$ , called the **growth function**, such that for any  $y \geq 0$  and  $\varepsilon > 0$ , if  $f_{\leq y} \neq \emptyset$ , then  $\lambda_f(y) \geq \text{diam}(f_{\leq y})/\zeta$ , where  $\zeta$  is an absolute constant, the **growth constant**, depending only on the family of functions and not on  $n$  and such that if  $\mathbf{q} \in \mathbb{R}^d$  with  $d(\mathbf{q}, f_{\leq y}) \leq \varepsilon \lambda_f(y)$ , then  $f(\mathbf{q}) \leq (1+\varepsilon)y$ . This is equivalent to  $f_{\leq y} \oplus \mathbf{B}(0, \varepsilon \lambda_f(y)) \subseteq f_{\leq (1+\varepsilon)y}$ , where  $\mathbf{B}(u, r)$  is the ball of radius  $r$  centered at  $u$ .
- (P3) **Existence of a sketch.** Given  $\delta > 0$  and a subset  $\mathcal{G} \subseteq \mathcal{F}$ , there is a  $\mathcal{H} \subseteq \mathcal{G}$  with  $|\mathcal{H}| = O(1/\delta^{\text{c}_{\text{sk}}})$  and  $y_0 = O(\text{cl}(\mathcal{G})(|\mathcal{G}|/\delta)^{\text{c}_{\text{sk}}})$  such that,  $\mathcal{H}$  is a  $(\delta, y_0)$ -sketch, where  $\text{c}_{\text{sk}}$  is some positive integer constant that depends on the given family of functions.

We also require some straightforward properties from the computation model:

- (C1)  $\forall \mathbf{q} \in \mathbb{R}^d$  and  $1 \leq i \leq n$ , the value  $f_i(\mathbf{q}) = d(\mathbf{q}, f_i)$  is computable in  $O(1)$  time.
- (C2) For any  $y \geq 0, r > 0$  and  $i$ , the set of grid cells approximating the sublevel set  $(f_i)_{\leq r}$  of  $f_i$ , that is  $(f_i)_{\leq y, \approx r} = \mathbf{G}_{\approx r}((f_i)_{\leq y})$  (see Definition II.1), is computable in linear time in its size.
- (C3) For any  $f_i, f_j \in \mathcal{F}, 1 \leq i, j \leq n$  the distance  $d(f_i, f_j)$  is computable in  $O(1)$  time.

We also assume that the growth function  $\lambda_{(f_i)}(y)$  from Condition (P2) be in fact computable easily, i.e., in  $O(1)$  time.

**Remark II.9.** We will use Condition (C2) for a given  $y$  and  $i$  only for  $r$  at least  $\Omega(\varepsilon \lambda_{(f_i)}(y))$ , i.e., we will use a grid on the sublevel set at a resolution typically  $\varepsilon$  times its growth function value at that level, which by Condition (C2) is also  $\Omega(\varepsilon \text{diam}((f_i)_{\leq y}))$ . As such the number of grid cells in the grid used is  $O(1/\varepsilon^d)$ .

**Remark II.10.** Condition (C2) when used with a resolution  $r = \varepsilon \lambda_{(f_i)}(y)$  requires that we precisely output the set of cells  $(f_i)_{\leq y, \approx r}$ . However, it is sufficient to output a set of cells at that resolution which contains  $(f_i)_{\leq y, \approx r}$ , is still contained in  $f_{\leq (1+\varepsilon)y}$  and is of size  $O(1/\varepsilon^d)$ . We state the stricter condition, because it is easier to digest. However it is easier and still correct to work with the relaxed condition.

### II-B1. Properties

The following are basic properties that the functions under consideration have. Their proofs are in the full version [1].

**Lemma II.11.** Let  $\mathcal{F}$  be a set of functions that satisfy the compactness (P1) and bounded growth (P2) conditions. Then, for any  $f \in \mathcal{F}$ , either  $f_{\leq 0} = \emptyset$  or  $f_{\leq 0}$  consists of a single point.

By the above lemma, we may assume that a symbolic perturbation guarantees that  $d(f, g) > 0$  for  $f \neq g$ . With this convention we have the following,

**Observation II.12.** If  $\text{cl}(\mathcal{G}) = 0$  for any non-empty subset  $\mathcal{G}$  then  $|\mathcal{G}| = 1$ .

We also assume that the quantities  $d(f, g)$  are distinct for all distinct pairs of functions.

**Lemma II.13.** Let  $f \in \mathcal{G}$  and  $y \geq 0$ . Suppose  $u, v \in f_{\leq y}$ . Then,  $uv \subseteq \mathcal{G}_{\leq (1+\zeta/2)y}$ , where  $uv$  denotes the segment joining  $u$  to  $v$ .

**Lemma II.14.** Let  $A_1, \dots, A_m \subseteq \mathbb{R}^d$  be compact connected sets. Let  $uv$  be any segment. Suppose that  $uv \cap A_i \neq \emptyset$  for all  $1 \leq i \leq k$  and  $uv \subseteq \bigcup_{i=1}^k A_i$ . Then, the sets  $A_i, 1 \leq i \leq k$ , are connected.

**Lemma II.15.** Suppose we are given  $\mathcal{H} \subseteq \mathcal{G} \subseteq \mathcal{F}$ ,  $\delta \geq 0$  and  $y \geq 0$ , and  $\mathcal{H}$  is a  $(\delta, y)$ -sketch for  $\mathcal{G}$ . Then,  $\text{cl}(\mathcal{H}) \leq (1+\delta)(1+\zeta/2) \max(y, \text{cl}(\mathcal{G}))$ .

The following testifies that a sketch approximates the distance to a set of functions.

**Lemma II.16.** Let  $\mathcal{H} \subseteq \mathcal{G}$  be sets of functions, where  $\mathcal{H}$  is a  $(\delta, y_0)$ -sketch for  $\mathcal{G}$  for some  $\delta \geq 0$  and  $y_0 \geq 0$ . Let  $\mathbf{q}$  be a point such that  $d(\mathbf{q}, \mathcal{G}) \geq y_0$ . Then we have that  $d(\mathbf{q}, \mathcal{H}) \leq (1+\delta)d(\mathbf{q}, \mathcal{G})$ .

### II-B2. Computing the connectivity level

We implicitly assume that the above relevant quantities can be computed efficiently. For example given some  $\delta > 0$ , and  $y_0$  as per the bound in condition (P3), a  $(\delta, y_0)$ -sketch can be computed in time  $O(|\mathcal{G}|/\delta^{\text{c}_{\text{sk}}})$  time. We also assume that  $\text{cl}(\mathcal{G})$  can be computed efficiently without resorting to the “brute force” method. The brute force method computes the individual distance of the functions and then computes a MST on the graph defined by vertices as the functions and edge lengths as their distance. Then  $\text{cl}(\mathcal{G})$  is the longest edge of this MST.

## III. CONSTRUCTING THE AVD

The input is a set  $\mathcal{F}$  of  $n$  functions satisfying the conditions of Section II-B, and a number  $0 < \varepsilon \leq 1$ . We preprocess  $\mathcal{F}$ , such that given a query point  $\mathbf{q}$  one can compute a  $f \in \mathcal{F}$ , where  $d(\mathbf{q}, \mathcal{F}) \leq d(\mathbf{q}, f) \leq (1+\varepsilon)d(\mathbf{q}, \mathcal{F})$ .

### III-A. Building blocks

#### III-A1. Near neighbor

Given a set of functions  $\mathcal{G}$ , a real number  $\alpha \geq 0$ , and a parameter  $\varepsilon > 0$ , a **near-neighbor** data-structure  $\mathcal{D}_{\text{near}} = \mathcal{D}_{\text{nr}}(\mathcal{G}, \varepsilon, \alpha)$  can decide (approximately) if a point has distance larger or smaller than  $\alpha$ . Formally, for a query point  $\mathbf{q}$  a near-neighbor query answers yes if  $d(\mathbf{q}, \mathcal{G}) \leq \alpha$ , and no if  $d(\mathbf{q}, \mathcal{G}) > (1+\varepsilon)\alpha$ . It can return either answer if  $d(\mathbf{q}, \mathcal{G}) \in \left( \alpha, (1+\varepsilon)\alpha \right]$ . If it returns yes, then it also returns a function  $f \in \mathcal{G}$  such that  $d(\mathbf{q}, f) \leq (1+\varepsilon)\alpha$ . The query time of this data-structure is denoted by  $T_{\leq}(m)$ , where  $m = |\mathcal{G}|$ . The proof of the following appears in the full version [1].

**Lemma III.1.** Given a set of  $m$  functions  $\mathcal{G} \subseteq \mathcal{F}$ ,  $\alpha > 0$  and  $\varepsilon > 0$ . One can construct a data structure (which is a compressed quadtree), of size  $O(m/\varepsilon^d)$ , in  $O(m\varepsilon^{-d} \log(m/\varepsilon))$  time, such that given any query point  $\mathbf{q} \in \mathbb{R}^d$  one can answer a  $(1+\varepsilon)$ -approximate near-neighbor query for the distance  $\alpha$ , in time  $T_{\leq}(m) = O(\log(m/\varepsilon))$ .

### III-A2. Interval data structure

Given a set of functions  $\mathcal{G}$ , real numbers  $0 < \alpha \leq \beta$ , and  $\varepsilon > 0$ , the **interval data structure** returns for a query point  $\mathbf{q}$ , one of the following:

- (A) If  $d(\mathbf{q}, \mathcal{G}) \in [\alpha, \beta]$ , then it returns a function  $g \in \mathcal{G}$  such that  $d(\mathbf{q}, g) \leq (1 + \varepsilon)d(\mathbf{q}, \mathcal{G})$ . It might also return such a function for values outside this interval.
- (B) “ $d(\mathbf{q}, \mathcal{G}) < \alpha$ ”. In this case it returns a function  $g \in \mathcal{G}$  such that  $d(\mathbf{q}, g) < \alpha$ .
- (C) “ $d(\mathbf{q}, \mathcal{G}) > \beta$ ”.

The time to perform an interval query is denoted by  $T_r(m, \alpha, \beta)$ . The proof of the following appears in the full version [1].

**Lemma III.2.** *Given a set of  $m$  functions  $\mathcal{G}$ , an interval  $[\alpha, \beta]$  and an approximation parameter  $\tau > 0$ , one can construct an interval data structure of size  $O(m\tau^{-d-1} \log(4\beta/\alpha))$ , in time  $O(m\tau^{-d-1} \log(4\beta/\alpha) \log(m/\tau))$ , such that given a query point  $\mathbf{q}$  one can answer  $(1 + \tau)$ -approximate nearest neighbor query for the distances in the interval  $[\alpha, \beta]$ , in time  $T_r(m, \alpha, \beta, f) = O\left(\log \frac{m \log(4\beta/\alpha)}{\tau}\right)$ .*

Lemma III.2 readily implies that if somehow a priori we know that the nearest neighbor distance lies in an interval of values of polynomial spread, then we would get the desired data-structure by just using Lemma III.2. To overcome this unbounded spread problem, we would first argue that, under our assumptions, there are only linear number of intervals where interesting things happen to the distance function.

### III-A3. Connected components of the sublevel sets

Given a finite set  $X$  and a partition of it into disjoint sets  $X = X_1 \cup \dots \cup X_k$ , let this partition be denoted by  $\langle X_1, \dots, X_k \rangle_X$ . For  $1 \leq i \leq k$ , each  $X_i$  is a **part** of the partition.

**Definition III.3.** *For two partitions  $P_A = \langle A_1, \dots, A_k \rangle_X$  and  $P_B = \langle B_1, \dots, B_l \rangle_X$  of the same set  $X$ ,  $P_B$  is a **refinement** of  $P_A$ , denoted by  $P_B \sqsubseteq P_A$ , if for any  $B_i$  there exists a set  $A_{j_i}$ , such that  $B_i \subseteq A_{j_i}$ . In the other direction,  $P_A$  is a **coarsening** of  $P_B$ .*

**Observation III.4.** *Given partitions  $\Pi, \Xi$  of a finite set  $X$ , if  $\Pi \sqsubseteq \Xi$  then  $|\Xi| \leq |\Pi|$ .*

**Definition III.5.** *Given partitions  $\Pi = \langle X_1, \dots, X_k \rangle_X \sqsubseteq \Xi = \langle X'_1, \dots, X'_{k'} \rangle_X$ , let  $\phi(\Pi, \Xi, i)$  be the function that return the set of indices of sets in  $\Pi$  whose union is  $X'_i \in \Xi$ .*

**Observation III.6.** *Given partitions  $\Pi \sqsubseteq \Xi$  of a set  $X$  with  $n$  elements. The partition function  $\phi(\Pi, \Xi, \cdot)$  can be computed in  $O(n)$  time. For any  $1 \leq i \leq |\Xi|$ , the set  $\phi(\Pi, \Xi, i)$  can be returned in  $O(|\phi(\Pi, \Xi, i)|)$  time, and its size can be returned in  $O(1)$  time.*

**Definition III.7.** *For  $\mathcal{G} \subseteq \mathcal{F}$  and  $\ell > 0$ , consider the intersection graph of the sets  $f_{\leq \ell}$ , for all  $f \in \mathcal{G}$ . Each connected component is a **cluster** of  $\mathcal{G}$  at level  $\ell$ . And the partition of  $\mathcal{G}$  by these clusters, denoted by  $C(\mathcal{G}, \ell)$ , is the  $\ell$ -*

**clustering** of  $\mathcal{G}$ .

The values  $\ell$  at which the  $\ell$ -clustering of  $\mathcal{F}$  changes are intuitively the critical values when the sublevel set of  $\mathcal{F}$  changes and which influence the AVD. These values are critical in trying to decompose the nearest neighbor search on  $\mathcal{F}$  into a search on smaller sets.

**Observation III.8.** *If  $0 \leq a \leq b$  then  $C(\mathcal{G}, a) \sqsubseteq C(\mathcal{G}, b)$ .*

The following lemma testifies that we can approximate the  $\ell$ -clustering quickly, for any number  $\ell$ . The proof is in the full version [1].

**Lemma III.9.** *Given  $\mathcal{G} \subseteq \mathcal{F}$ ,  $\ell \geq 0$ , and  $\varepsilon > 0$ , one can compute, in  $O\left(\frac{m}{\varepsilon^d} \log(m/\varepsilon)\right)$  time, a partition  $\Psi = \Psi_\varepsilon(\mathcal{G}, \ell)$ , such that  $C(\mathcal{G}, \ell) \sqsubseteq \Psi \sqsubseteq C(\mathcal{G}, (1 + \varepsilon)\ell)$ , where  $m = |\mathcal{G}|$ .*

**Remark III.10.** *The partition  $\Psi$  computed by Lemma III.9 is monotone, that is, for  $\ell \leq \ell'$  and  $\varepsilon \leq \varepsilon'$ , we have  $\Psi_\varepsilon(\mathcal{G}, \ell) \sqsubseteq \Psi_{\varepsilon'}(\mathcal{G}, \ell')$ . Moreover, for each cluster  $C \in \Psi_\varepsilon(\mathcal{G}, \ell)$ , we have that  $\text{cl}(C) \leq (1 + \varepsilon)\ell$ .*

### III-A4. Computing a splitting distance

**Definition III.11.** *Given a partition  $\Psi = \Psi_\varepsilon(\mathcal{G}, \ell)$  of  $\mathcal{G}$ , with  $m = |\Psi|$  clusters, a distance  $x$  is a **splitting distance** if  $m/4 \leq |\Psi_1(\mathcal{G}, x/4)|$  and  $|\Psi_1(\mathcal{G}, x)| \leq (7/8)m$ .*

**Lemma III.12.** *Given a partition  $\Psi = \Psi_\varepsilon(\mathcal{G}, \ell)$  of  $\mathcal{G}$ , one can compute a splitting distance for it, in expected  $O(n(\log n + t))$  time, where  $n = |\mathcal{G}|$  and  $t$  is the maximum cluster size in  $\Psi$ .*

*Proof:* For each cluster  $C \in \Psi$ , let  $r_C$  be its distance from all the functions in  $\mathcal{G} \setminus C$ ; that is  $r_C = \min_{f \in C} \min_{g \in \mathcal{G} \setminus C} d(f, g)$ . Note that  $r_C \geq \ell$ . Now, let  $r_1 \leq r_2 \leq \dots \leq r_m$  be these distances for the  $m$  clusters of  $\Psi$ . We randomly pick a cluster  $C \in \Psi$  and compute  $\ell' = r_C$  for it by brute force – computing the distance of each function of  $C$  with the functions of  $\mathcal{G} \setminus C$ .

Let  $i$  be the rank of  $\ell' = r_C$  among  $r_1, \dots, r_m$ . With probability 1/2, we have that  $m/4 \leq i \leq (3/4)m$ . If so we have that:

- (A) All the clusters that correspond to  $r_i, \dots, r_m$  are singletons in the partition  $\Psi_1(\mathcal{G}, \ell'/4)$ , as the distance of each one of these clusters is larger than  $\ell'$ . We conclude that  $|\Psi_1(\mathcal{G}, \ell'/4)| \geq m/4$ .
- (B) All the clusters of  $\Psi$  that correspond to  $r_1, \dots, r_i$  are contained inside a larger cluster of  $\Psi_1(\mathcal{G}, \ell')$  (i.e., they were merged with some other cluster). But then, the number of clusters in  $\Psi_1(\mathcal{G}, \ell')$  is at most  $(7/8)m$ . Indeed, put an edge between such a cluster, to the cluster realizing the smallest distance with it. This graph has at least  $e \geq m/4$  edges, and it is easy to see that each component of size at least 2 in the underlying undirected graph has the same number of edges as vertices. As such the number of singleton components is at most  $m - e$  while the number of components of size at least 2 is at most  $e/2$ . It follows that the total number of components is at most  $m - e/2 \leq 7m/8$ . Since each such component corresponds to a cluster in  $\Psi_1(\mathcal{G}, \ell')$  the claim is

```

Search( $\mathcal{G}, \Upsilon, \mathbf{q}$ )
//  $\mathcal{G}$ : set of functions
//  $\Upsilon = \Psi_1(\mathcal{G}, \ell)$  for some value  $\ell$ 
// Invariant:  $d(\mathbf{q}, \mathcal{G}) > \ell N$ 
if  $|\Upsilon| = 1$  then
    return  $d(\mathbf{q}, \mathcal{G}) = \min_{f \in \mathcal{G}} d(\mathbf{q}, f)$  (*)
 $x \leftarrow$  compute a splitting distance of  $\Upsilon$ , see Lemma III.12

// Perform an interval approximate nearest
// neighbor query on the interval  $[x/8N, x8N]$ 
// for the set  $\mathcal{G}$ , see Lemma III.2.
if  $d(\mathbf{q}, \mathcal{G}) \in [x/8N, x8N^2]$  or  $(1 + \frac{\varepsilon}{4})$ -ANN found then
    return nearest function found by the
         $(1 + \varepsilon/4)$ -approximate interval query.
if  $d(\mathbf{q}, \mathcal{G}) < x/8N$  then
     $f \leftarrow$  2-approximate near neighbor query on  $\mathcal{G}$ 
        and distance  $x/8$ , see Lemma III.1.
    Find cluster  $C \in \Psi_1(\mathcal{G}, x/4)$ , such that  $f \in C$ ,
        see Lemma III.9.
    return Search( $C, \Upsilon[C], \mathbf{q}$ )
if  $d(\mathbf{q}, \mathcal{G}) > x8N^2$  then
    return Search( $\mathcal{G}, \Psi_1(\mathcal{G}, xN), \mathbf{q}$ ) (**)

```

Figure III.1: Search algorithm: We are given a query point  $\mathbf{q}$ , and an approximation parameter  $\varepsilon > 0$ . The quantity  $N$  is a parameter to be specified shortly. Initially, we call this procedure on the set of functions  $\mathcal{F}$  with  $\Upsilon$  being the partition of  $\mathcal{F}$  into singletons (i.e.,  $\ell = 0$ ). Here,  $\Upsilon[C]$  denotes the partition of  $C$  induced by the partition  $\Upsilon$ .

proved.

Now, compute  $\Psi_1(\mathcal{G}, \ell')$  and  $\Psi_1(\mathcal{G}, \ell'/4)$  using Lemma III.9. With probability at least half they have the desired sizes, and we are done. Otherwise, we repeat the process. In each iteration we spend  $O(n(\log n + t))$  time, and the probability for success is half. As such, in expectation the number of rounds needed is constant. ■

### III-B. The search procedure

#### III-B1. An initial “naive” implementation

The search procedure is presented in Figure III.1. The following lemma, whose proof appears in the full version [1], proves its correctness.

**Lemma III.13.** **Search**( $\mathcal{G}, \Upsilon, \mathbf{q}$ ) returns a function  $f \in \mathcal{G}$ , such that  $d(\mathbf{q}, f) \leq (1 + \varepsilon)d(\mathbf{q}, \mathcal{G})$ . The depth of the recursion of **Search** is  $h = O(\log n)$ , where  $n = |\mathcal{G}|$ .

#### III-B2. But where is the beef? Modifying **Search** to provide fast query time

The reader might wonder how we are going to get an efficient search algorithm out of **Search**, as the case that  $\Upsilon$  is a single cluster, still requires us to perform a scan on all the functions in this cluster and compute their distance from the query point  $\mathbf{q}$ . Note however, we have the invariant that the distance of interest is polynomially larger than the connectivity level of each of the clusters of  $\Upsilon$ . In particular, precomputing for all the sets of functions such that (\*) might

be called on, their  $\varepsilon/8$ -sketches, and answering the query by computing the distance on the sketches, reduces the query time to  $O(1/\varepsilon^{c_{\text{sk}}} + \log^2 n)$  (assuming that we precomputed all the data-structures used by the query process). Indeed, an interval query takes  $O(\log n)$  time, and there  $O(\log n)$  such queries. The final query on the sketch takes time proportional to the sketch size which is  $O(1/\varepsilon^{c_{\text{sk}}})$ .

As such, the major challenge is not making the query process fast, but rather building the search structure quickly, and arguing that it requires little space.

#### III-B3. Sketching a sketch

To improve the efficiency of the preprocessing for **Search**, we are going to use sketches more aggressively. Specifically, for each of the clusters of  $\Upsilon$ , we can compute their  $\delta$ -sketches, for  $\delta = \varepsilon/(8h) = O(\varepsilon/\log n)$ , see Lemma III.13. From this point on, when we manipulate this cluster, we do it on its sketch. To make this work set  $N = n^{4c_{\text{sk}}}$ , see (P3)<sub>p4</sub> and Lemma II.15.

The only place in the algorithm where we need to compute the sketches, is in (\*\*) in Figure III.1. Specifically, we compute  $\Psi_1(\mathcal{G}, xN)$ , and for each new cluster  $C \in \Psi_1(\mathcal{G}, xN)$ , we combine all the sketches of the clusters  $D \in \Upsilon$  such that  $D \subseteq C$  into a single set of functions. We then compute a  $\delta$ -sketch for this set, and this sketch is this cluster from this point on. In particular, the recursive calls to **Search** would send the sketches of the clusters, and not the clusters themselves. Conceptually, the recursive call would also pass the minimum distance where the sketches are active – it is easy to verify that we use these sketches only at distances that are far away and are thus allowable (i.e., the sketches represent the functions they correspond to, well in these distances).

Importantly, whenever we compute such a new set, we do so for a distance that is bigger by a polynomial factor (i.e.,  $N$ ) than the values used to create the sketches of the clusters being merged. Indeed, observe that  $x > \ell$  and as such  $xN$  is  $N$  times bigger than  $\ell$  (an upper bound on the value used to compute the input sketches).

As such, all these sketches are valid, and can be used at this distance (or any larger distance). Of course, the quality of the sketch deteriorates. In particular, since the depth of recursion is  $h$ , the worst quality of any of the sketches created in this process is at most  $(1 + \delta)^h \leq 1 + \varepsilon/4$ .

Significantly, before using such a sketch, we would shrink it by computing a  $\varepsilon/8$ -sketch of it. This would reduce the sketch size to  $O(1/\varepsilon^{c_{\text{sk}}})$ . Note, however, that this still does not help us as far as recursion – we must pass the larger  $\delta$ -sketches in the recursive call of (\*\*).

This completes the description of the search procedure. It is still unclear how to precompute all the data-structures required during the search. To do that, we need to better understand what the search process does.

### III-C. The connectivity tree, and the preprocessing

Given a set of functions  $\mathcal{F}$ , consider the tree tracking the connected components of the MST of the functions. Formally, initially we start with  $n$  singletons (which are the leaves of the tree) that are labeled with the value zero, and we store them

in a set  $\mathcal{F}$  of active nodes. Now, we compute for each pair of sets of functions  $X, Y \in \mathcal{F}$  the distance  $d(X, Y)$ , and let  $X', Y'$  be the pair realizing the minimum of this quantity. Merge the two sets into a new set  $Z = X' \cup Y'$ , create a new node for this set having the node for  $X'$  and  $Y'$  as children, and set its label to be  $d(X', Y')$ . Finally, remove  $X'$  and  $Y'$  from  $\mathcal{F}$  and insert  $Z$  into it. Repeat till there is a single element in  $\mathcal{F}$ . Clearly, the result is a tree that tracks the connected components of the MST.

To make the presentation consistent, let  $d_{\approx}(X, Y)$  be the minimum  $x$  such that  $\Psi_1(X \cup Y, x)$  is connected. Computing  $d_{\approx}(X, Y)$  can be done by computing  $d_{\approx}(f, g)$  for each pair of functions separately. This in turn, can be done by first computing  $\alpha = d(f, g)$  and observing that  $r$  is between  $\alpha/4$  and  $\alpha$ . In particular,  $r$  must be a power of two, so there are only 3 candidate values to consider, and which is the right one can be decided using Lemma III.9.

So, in the above, we use  $d_{\approx}(\cdot, \cdot)$  instead of  $d(\cdot, \cdot)$ , and let  $\mathcal{H}$  be the resulting tree. For a value  $\ell$ , let  $L_{\mathcal{H}}(\ell)$  be the set of nodes such that their label is smaller than  $\ell$ , but their parent label is larger than  $\ell$ . It is easy to verify that  $L_{\mathcal{H}}(\ell)$  corresponds to  $\Psi = \Psi_1(\mathcal{F}, \ell)$ ; indeed, every cluster  $C \in \Psi$  corresponds to a node  $u \in L_{\mathcal{H}}(\ell)$ , such that the set of functions stored in the leaves of the subtree of  $u$ , denoted by  $F(u)$  is  $C$ . The following can be easily proved by induction.

**Lemma III.14.** *Consider a recursive call **Search**( $\mathcal{G}, \Upsilon, \mathbf{q}$ ) made during the search algorithm execution. Then  $\mathcal{G} = F(u)$ , and  $\Upsilon = \left\{ F(v) \mid v \in L_{\mathcal{H}}(\ell) \text{ and } v \text{ is in the subtree of } u \right\}$ .*

*That is, a recursive call of **Search** corresponds to a subtree of  $\mathcal{H}$ .*

Of course, not all possible subtrees are candidates to be such a recursive call. In particular, **Search** can now be interpreted as working on a subtree  $T$  of  $\mathcal{H}$ , as follows:

- (A) If  $T$  is a single node  $u$ , then find the closet function to  $F(u)$ . Using the sketch this can be done quickly.
- (B) Otherwise, compute a distance  $x$ , such that the number of nodes in the level  $L_T(x)$  is roughly half the number of leaves of  $T$ .
- (C) Using interval data-structure determine if the distance  $d(\mathbf{q}, F(T))$  is in the range  $[x/8N, x8N^2]$ . If so, we found the desired ANN.
- (D) If  $d(\mathbf{q}, F(T)) > x8N^2$  then continue recursively on portion of  $T$  above  $L_T(x)$ .
- (E) If  $d(\mathbf{q}, F(T)) < x/8N$  then we know the node  $u \in L_T(x)$  such that the ANN belongs to  $F(u)$ . Continue the search recursively on the subtree of  $T$  rooted at  $u$ .

That is, **Search** breaks  $T$  into subtrees, and continues the search recursively on one of the subtrees. Significantly, every such subtree has constant fraction of the size of  $T$ , and every edge of  $T$  belongs to a single such subtree.

The preprocessing now works by precomputing all the data-structures required by **Search**. Of course, the most natural approach would be to precompute  $\mathcal{H}$ , and build the search tree by simulating the above recursion on  $\mathcal{H}$ . Fortunately, this is not necessary, we simulate running **Search**,

and investigate all the different recursive calls. We thus only use the above  $\mathcal{H}$  in analyzing the preprocessing running time. See Figure III.1p6.

In particular, given a subtree  $T$  with  $m$  edges, the corresponding partition  $\Upsilon$  would have at most  $m$  sets. Each such set would have a  $\delta$ -sketch, and we compute a  $\varepsilon/8$ -sketch for each one of these sketches. Namely, the input size here is  $M = O(m/\delta^{c_{\text{sk}}})$ . Computing the  $\varepsilon/8$ -sketches for each one of these sketches reduces the total number of functions to  $M' = O(m/\varepsilon^{c_{\text{sk}}})$ , and takes  $U_1 = O(M/\varepsilon^{c_{\text{sk}}}) = O(m(\varepsilon\delta)^{-c_{\text{sk}}})$  time, see Section II-B2. Computing the splitting distance, using Lemma III.12, takes  $U_2 = O(M' \log M' + 1/\varepsilon^{c_{\text{sk}}}) = O(m\varepsilon^{-c_{\text{sk}}} \log m)$  time. Computing the interval data-structure Lemma III.2 takes  $U_3 = O(M'\varepsilon^{-d-1} \log n \log M')$  time, and requires  $S_1 = O(M'\varepsilon^{-d-1} \log n)$  space. This breaks  $T$  into edge disjoint subtrees  $T_1, \dots, T_t$ , and we compute the search data-structure for each one of them separately (each one of these subtrees is smaller by a constant fraction of the original tree). Finally, we need to compute the  $\delta$ -sketches for the clusters sent to the appropriate recursive calls, and this takes  $U_4 = O(M/\delta^{c_{\text{sk}}})$ , by Section II-B2.

Every edge of the tree  $T$  gets charged for the amount of work spent in building the top level data-structure. That is, the top level amortized work each edge of  $T$  has to pay is

$$\begin{aligned} & O\left((U_1 + U_2 + U_3 + U_4)/m\right) \\ &= O\left((\varepsilon\delta)^{-c_{\text{sk}}} + \varepsilon^{-c_{\text{sk}}} \log m + \varepsilon^{-d-1-c_{\text{sk}}} \log^2 n + \delta^{-2c_{\text{sk}}}\right) \\ &= O\left(\varepsilon^{-2c_{\text{sk}}} \log^{2c_{\text{sk}}} n\right), \end{aligned}$$

assuming  $c_{\text{sk}} \geq 2$ . Since an edge of  $T$  gets charged at most  $O(\log n)$  times by this recursive construction, we conclude that the total preprocessing time is  $O(n\varepsilon^{-2c_{\text{sk}}} \log^{2c_{\text{sk}}+1} n)$ .

As for the space, we have by the same argumentation, that each edge requires  $O(\log n \cdot (S_1/m)) = (\varepsilon^{-d-1-c_{\text{sk}}} \log^2 n)$ . As such, the overall space used by the data-structure is  $(n\varepsilon^{-d-1-c_{\text{sk}}} \log^2 n)$ . As for the query time, it boils down to  $O(\log n)$  interval queries, and then scanning one  $O(\varepsilon)$ -sketch. As such, this takes  $O(\log^2 n + 1/\varepsilon^{c_{\text{sk}}})$  time.

### III-D. The result

Our main result is the following,

**Theorem III.15.** *Let  $\mathcal{F}$  be a set of  $n$  functions in  $\mathbb{R}^d$  that complies with our assumptions, see Section II-B, and has sketch constant  $c_{\text{sk}} \geq d$ . Then, one can build a data-structure to answer ANN for this set of functions, with the following properties:*

- (A) *The query time is  $O(\log n + 1/\varepsilon^{c_{\text{sk}}})$ .*
- (B) *The preprocessing time is  $O(n\varepsilon^{-2c_{\text{sk}}} \log^{2c_{\text{sk}}+1} n)$ .*
- (C) *The space used is  $O(n\varepsilon^{-d-1-c_{\text{sk}}} \log^2 n)$ .*

*Proof:* The query time for the algorithm described above is  $O(\log^2 n + 1/\varepsilon^{c_{\text{sk}}})$ . To get the improved query time, we observe that **Search** performs a sequence of point-location queries in a sequence of interval near neighbor data-structures (i.e., compressed quadtrees), and then it scans a set of functions of size  $O(1/\varepsilon^{c_{\text{sk}}})$  to find the ANN. We take all these quadtrees spread through our data-structure, and assign

them priority, where a quadtree  $\mathcal{Q}_1$  has higher priority than a compressed quadtree  $\mathcal{Q}_2$  if  $\mathcal{Q}_1$  is queried after  $\mathcal{Q}_2$  for any search query. This defines an acyclic ordering on these compressed quadtrees. Overlaying all these compressed quadtrees together, one needs to return for the query point, the leaf of the highest priority quadtree that contains the query point. This can be easily done by scanning the compressed quadtree, and for every leaf computing the highest priority leaf that contains it (observe, that here we are overlaying only the nodes in the compressed quadtrees that are marked by some sublevel set – nodes that are empty are ignored).

A tedious but straightforward induction implies that doing a point-location query in the resulting quadtree is equivalent to running the search procedure as described above. Once we found the leaf that contains the query point, we scan the sketch associated with this cell, and return the computed nearest-neighbor. ■

An important corollary is the following,

**Corollary III.16.** *Let  $\mathcal{F}$  be a set of  $n$  functions in  $\mathbb{R}^d$  that complies with our assumptions, see Section II-B, and has sketch constant  $c_{\text{sk}} \geq d$ . Then, one can build a data-structure to answer ANN for this set of functions, with the following properties:*

- (A) *The improved query time is  $O(\log n)$ .*
- (B) *The preprocessing time is  $O(n/\varepsilon^{O(1)} \log^{2c_{\text{sk}}+1} n)$ .*
- (C) *The space used is  $S = O(n/\varepsilon^{O(1)} \log^2 n)$ .*

*In particular, we can compute an AVD of complexity  $O(S)$  for the given functions. That is, one can compute a space decomposition, such that every region has a single function associated with it, and for any point in this region, this function is a  $(1 + \varepsilon)$ -ANN among the functions of  $\mathcal{F}$ . Here, a region is either a cube, or the set difference of two cubes.*

*Proof:* We build the data-structure of Theorem III.15, except that instead of linearly scanning the sketch during the query time, we preprocess each such sketch for an exact point-location query; that is, we compute the lower envelope of the sketch and preprocess it for vertical ray shooting [25]. This would require  $O(1/\varepsilon^{O(1)})$  space for each such sketch, and the linear scanning that takes  $O(1/\varepsilon^{O(1)})$  time, now is replaced by a point-location query that takes  $O(\log 1/\varepsilon^{O(1)}) = O(\log 1/\varepsilon) = O(\log n)$ , as desired.

As for the second part, observe that every leaf of the compressed quadtree is the set difference of two canonical grid cells. The lower envelope of the functions associated with such a leaf, induce a partition of this leaf into regions with total complexity  $O(1/\varepsilon^{O(1)})$ . ■

#### IV. APPLICATIONS

We present some concrete classes of functions that satisfy our framework, and for which we construct AVD's efficiently.

##### IV-A. Multiplicative distance functions with additive offsets

As a warm-up we present the simpler case of additively offset multiplicative distance functions. The results of this section are almost subsumed by more general results in Section IV-B. Here the sublevel sets look like expanding

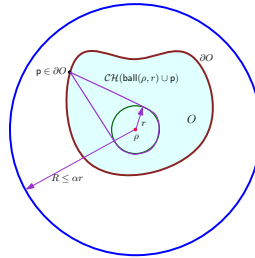


Figure IV.1: Being  $\alpha$ -rounded fat.

balls but there is a time lag before the balls even come into existence i.e. sublevel sets are empty up-to a certain level, this corresponds to the additive offsets. In Section IV-B the sublevel sets are more general fat bodies but there is no additive offset. The results in the present section essentially give an AVD construction of approximate weighted Voronoi diagrams. More formally, we are given a set of points  $P = \{p_1, \dots, p_n\}$ . For  $i = 1, \dots, n$ , the point  $p_i$  has weight  $w_i > 0$ , and a constant  $\alpha_i \geq 0$  associated with it. We define  $f_i(q) = w_i \|q - p_i\| + \alpha_i$ . Let  $\mathcal{F} = \{f_1, \dots, f_n\}$ . We have,  $(f_i)_{\leq y} = \emptyset$  for  $y < \alpha_i$  and  $(f_i)_{\leq y} = B(p_i, \frac{y - \alpha_i}{w_i})$  for  $y \geq \alpha_i$ .

The proof that  $\mathcal{F}$  complies with the framework of Section II-B is present in the full version [1].

We thus get the following result.

**Theorem IV.1.** *Consider a set  $P$  of  $n$  points in  $\mathbb{R}^d$ , where the  $i$ th point  $p_i$  has additive weight  $\alpha_i \geq 0$  and multiplicative weight  $w_i > 0$ . The  $i$ th point induces the additive/multiplicative distance function  $f_i(q) = w_i \|q - p_i\| + \alpha_i$ . Then one can compute a  $(1 + \varepsilon)$ -AVD for these distance functions, with near linear space complexity, and logarithmic query time. See Theorem III.15 for the exact bounds.*

##### IV-B. Scaling distance – generalized polytope distances

Let  $O \subseteq \mathbb{R}^d$  be a compact set homeomorphic to  $B(0, 1)$  and containing a “center” point  $\rho$  in its interior. Then  $O$  is **star shaped** if for any point  $v \in O$  the entire segment  $\rho v$  is also in  $O$ . Naturally, any convex body  $O$  with any center  $\rho \in O$  is star shaped. The  **$t$ -scaling** of  $O$  with a center  $\rho$  is the set  $tO = \{t(v - \rho) + \rho \mid v \in O\}$ .

Given a star shaped object  $O$  with a center  $\rho$ , the **scaling distance** of a point  $q$  from  $O$  is the minimum  $t$ , such that  $p \in tO$ , and let  $F_O(q)$  denote this distance function. Note that, for any  $y \geq 0$ , the sublevel set  $(F_O)_{\leq y}$  is the  $y$ -scaling of  $O$ , that is  $(F_O)_{\leq y} = yO$ .

Note that for a point  $p \in \mathbb{R}^d$ , if we take  $O = B(p, 1)$  with center  $p$ , then  $F_O(q) = \|p - q\|$ . That is, this distance notion is a strict extension of the Euclidean distance.

Henceforward, for this section, we assume that an object  $O$  contains the origin in its interior and the origin is the designated center, unless otherwise stated.

**Definition IV.2.** *Let  $O \subseteq \mathbb{R}^d$  be a star shaped object centered at  $p$ . We say that  $O$  is  $\alpha$ -fat if there is a number  $r$  such that,  $\text{ball}(p, r) \subseteq O \subseteq \text{ball}(p, \alpha r)$ .*



**Definition IV.3.** Let  $O$  be a star shaped object centered at  $\rho$ . We say that  $O$  is  $\alpha$ -rounded fat if there is a radius  $r$  such that, (i)  $\text{ball}(\rho, r) \subseteq O \subseteq \text{ball}(\rho, \alpha r)$  and, (ii) For every point  $p$  in the boundary of  $O$ , the cone  $\mathcal{CH}(\text{ball}(\rho, r) \cup p)$ , lies within  $O$ , see Figure IV.1.

By definition any  $\alpha$ -rounded fat object is also  $\alpha$ -fat. However, it is not true that a  $\alpha$ -fat object is necessarily  $\alpha'$ -rounded fat for any  $\alpha'$ , that is even allowed to depend on  $\alpha$ . The following useful result is easy to see.

**Lemma IV.4.** Let  $O$  be a  $\alpha$ -fat object. If  $O$  is convex then  $O$  is also  $\alpha$ -rounded fat.

Given a set  $\mathcal{O} = \{O_1, O_2, \dots, O_n\}$  of  $n$  star shaped objects, consider the set  $\mathcal{F}$  of  $n$  scaling distance functions, where the  $i$ th function, for  $i = 1, \dots, n$  is  $f_i = F_{O_i}$ . We assume that the boundary of each object  $O_i$  has constant complexity.

The proof that  $\mathcal{F}$  complies with the framework of Section II-B is present in the full version [1].

We conclude that for  $\alpha$ -rounded fat objects, the scaling distance functions they define, falls under our framework. We thus get the following result.

**Theorem IV.5.** Consider a set  $\mathcal{O}$  of  $\alpha$ -rounded fat objects in  $\mathbb{R}^d$ , for some constant  $\alpha$ . Then one can compute a  $(1 + \varepsilon)$ -AVD for the scaling distance functions induced by  $\mathcal{O}$ , with near linear space complexity, and logarithmic query time. See Theorem III.15 and Corollary III.16 for the exact bounds.

Note, that the result in Theorem IV.5 covers any symmetric convex metric. Indeed, given a convex symmetric shape  $C$  centered at the origin, the distance it induces for any pair of points  $\mathbf{p}, \mathbf{u} \in \mathbb{R}^d$ , is the scaling distance of  $C$  centered  $\mathbf{p}$  to  $\mathbf{u}$  (or, by symmetry, the scaling distance of  $\mathbf{p}$  from  $C$  centered at  $\mathbf{u}$ ). Under this distance  $\mathbb{R}^d$  is a metric space, and of course, the triangle inequality holds. By an appropriate scaling of space, which does not affect the norm (except for scaling it) we can make  $C$  fat, and now Theorem IV.5 applies. Of course, Theorem IV.5 is considerably more general, allowing each of the points to induce a different scaling distance function, and the distance induced does not have to obey the triangle inequality.

#### IV-C. Nearest furthest-neighbor

For a set of points  $S \subseteq \mathbb{R}^d$  and a point  $\mathbf{q}$ , the **furthest-neighbor distance** of  $\mathbf{q}$  from  $S$ , is  $\mathcal{F}_S(\mathbf{q}) = \max_{s \in S} \|\mathbf{q} - \mathbf{s}\|$ ; that is, it is the furthest one might have to travel from  $\mathbf{q}$  to arrive to a point of  $S$ . For example,  $S$  might be the set of locations of facilities, where it is known that one of them is always open, and one is interested in the worst case distance a client has to travel to reach an open facility. The function  $\mathcal{F}_S(\cdot)$  is known as the **furthest-neighbor Voronoi** diagram, and while its worst case combinatorial complexity is similar to the regular Voronoi diagram, it can be approximated using a constant size representation (in low dimensions), see [26].

Given  $n$  sets of points  $P_1, \dots, P_n$  in  $\mathbb{R}^d$ , we are interested in the distance function  $\mathcal{F}(\mathbf{q}) = \min_i \mathcal{F}_i(\mathbf{q})$ , where  $\mathcal{F}_i(\mathbf{q}) = \mathcal{F}_{P_i}(\mathbf{q})$ . This quantity arises naturally when one tries to model uncertainty; indeed, let  $P_i$  be the set of possible locations

of the  $i$ th point (i.e., the location of the  $i$ th point is chosen randomly, somehow, from the set  $P_i$ ). Thus,  $\mathcal{F}_i(\mathbf{q})$  is the worst case distance to the  $i$ th point, and  $\mathcal{F}(\mathbf{q})$  is the worst-case nearest neighbor distance to the random point-set generated by picking the  $i$ th point from  $P_i$ , for  $i = 1, \dots, n$ . We refer to  $\mathcal{F}(\cdot)$  as the **nearest furthest-neighbor** distance, and we are interested in its approximation.

A naive solution to this problem would maintain a data structure for computing the furthest neighbor approximately for each of the  $P_i$  and then just compute the minimum of those distances. A data-structure to compute a  $1 - \varepsilon$  approximation to the furthest neighbor takes  $O(1/\varepsilon^d)$  space for  $O(1/\varepsilon^d)$  query time, see [26] although this was probably known before. Thus the entire data structure would take up total space of  $O(n/\varepsilon^d)$  with a query time of  $O(n/\varepsilon^d)$ . By using our general framework we can speed up the computation. We can show that  $\mathcal{F}_i$ , for  $i = 1, \dots, n$  satisfy the conditions (P1) – (P3) and (C1)–(C3). The details are present in the full version [1].

We thus get the following result.

**Theorem IV.6.** Given  $n$  point sets  $P_1, \dots, P_n$  in  $\mathbb{R}^d$  with a total of  $m$  points, and a parameter  $\varepsilon > 0$ , one can preprocess the points into an AVD, of size  $\tilde{O}(n)$ , for the nearest furthest-neighbor distance defined by these point sets. One can now answer  $(1 + \varepsilon)$ -approximate NN queries for this distance in  $O(\log n)$  time. (Note, that the space and query time used, depend only on  $n$ , and not on the input size.)

*Proof:* We only need to show how to get the improved space and query time. Observe that every one of the sets  $P_i$  can be replaced by a subset  $S_i \subseteq P_i$ , of size  $O(1/\varepsilon^d \log(1/\varepsilon))$ , such that for any point  $\mathbf{q} \in \mathbb{R}^d$ , we have that  $\mathcal{F}_{S_i}(\mathbf{q}) \leq \mathcal{F}_{P_i}(\mathbf{q}) \leq (1 + \varepsilon/4)\mathcal{F}_{S_i}(\mathbf{q})$ . Such a subset can be computed in  $O(|P_i|)$  time, see [26]<sup>1</sup>. We thus perform this transformation for each one of the uncertain point sets  $P_1, \dots, P_n$ , which reduces the input size to  $O(n/\varepsilon^d \log(1/\varepsilon))$ . We now apply our main result to the distance functions induced by the reduced sets  $S_1, \dots, S_n$ . ■

## V. CONCLUSIONS

In this paper, we investigated what classes of functions have minimization diagrams that can be approximated efficiently – where our emphasis was on distance functions. We defined a general framework and the requirements on the distance functions to fall under it. For this framework, we presented a new data-structure, with near linear space and preprocessing time. This data-structure can evaluate (approximately) the minimization diagram of a query point in logarithmic time. Surprisingly, one gets an AVD (approximate Voronoi diagram) of this complexity; that is, a decomposition of space with near linear complexity, such that for every region of this decomposition a single function serves as an ANN for all points in this region.

We also showed some interesting classes of functions for which we get this AVD. For example, additive and multi-

<sup>1</sup>One computes an appropriate exponential grid, of size  $O(1/\varepsilon^d \log(1/\varepsilon))$ , and picks from each grid cell one representative point from the points stored inside this cell.

plicative weighted distance functions. No previous results of this kind were known, and even in the plane, multiplicative Voronoi diagrams have quadratic complexity in the worst case (for which the AVD generated has near linear complexity for any constant dimension). The framework also works for Minkowski metrics of fat convex bodies, and nearest furthest-neighbor. However, our main result applies to even more general distance functions.

Several questions remain open for further research:

- (A) Are the additional polylog factors in the space necessary? In particular, it seems unlikely that using WSPD's directly, as done by Arya and Malamatos [17], should work in the most general settings, so reducing the logarithmic dependency seems quite interesting. Specifically, can the Arya and Malamatos construction [17] be somehow adapted to this framework, possibly with some additional constraints on the functions, to get a linear space construction?
- (B) On the applications side, are constant degree polynomials a good family amenable to our framework? Specifically, consider a polynomial  $\tau(x)$  that is positive for all  $x \geq 0$ . Given a point  $u$ , we associate the distance function  $f(q) = \tau(\|q - u\|)$  with  $u$ . Given a set of such distance functions, under which conditions, can one build an AVD for these functions efficiently? (It is not hard to see that in the general case this is not possible, at least under our framework.)

#### REFERENCES

- [1] S. Har-Peled and N. Kumar, "Approximating minimization diagrams and generalized proximity search," *CoRR*, 2013, <http://arxiv.org/abs/1304.0393>.
- [2] J. Matoušek, "Efficient partition trees," *Discrete Comput. Geom.*, vol. 8, pp. 315–334, 1992.
- [3] T. M. Chan, "Optimal partition trees," in *Proc. 26th Annu. ACM Sympos. Comput. Geom.*, 2010, pp. 1–10.
- [4] M. Sharir and P. K. Agarwal, *Davenport-Schinzel Sequences and Their Geometric Applications*. New York: Cambridge University Press, 1995. [Online]. Available: <http://us.cambridge.org/titles/catalogue.asp?isbn=0521470250>
- [5] M. de Berg, O. Cheong, M. van Kreveld, and M. H. Overmars, *Computational Geometry: Algorithms and Applications*, 3rd ed. Springer-Verlag, 2008. [Online]. Available: <http://www.cs.uu.nl/geobook/>
- [6] K. L. Clarkson, "A randomized algorithm for closest-point queries," *SIAM J. Comput.*, vol. 17, pp. 830–847, 1988.
- [7] P. K. Agarwal and J. Matoušek, "Ray shooting and parametric search," *SIAM J. Comput.*, vol. 22, pp. 540–570, 1993.
- [8] S. Meiser, "Point location in arrangements of hyperplanes," *Inform. Comput.*, vol. 106, pp. 286–303, 1993.
- [9] K. L. Clarkson, "Nearest-neighbor searching and metric space dimensions," in *Nearest-Neighbor Methods for Learning and Vision: Theory and Practice*, G. Shakhnarovich, T. Darrell, and P. Indyk, Eds. MIT Press, 2006, pp. 15–59.
- [10] P. Indyk and R. Motwani, "Approximate nearest neighbors: Towards removing the curse of dimensionality," in *Proc. 30th Annu. ACM Sympos. Theory Comput.*, 1998, pp. 604–613. [Online]. Available: <http://theory.lcs.mit.edu/indyk/nndraft.ps>
- [11] S. Har-Peled, P. Indyk, and R. Motwani, "Approximate nearest neighbors: Towards removing the curse of dimensionality," *Theory Comput.*, vol. 8, pp. 321–350, 2012, special issue in honor of Rajeev Motwani.
- [12] A. Andoni and P. Indyk, "Near-optimal hashing algorithms for approximate nearest neighbor in high dimensions," *Commun. ACM*, vol. 51, no. 1, pp. 117–122, 2008.
- [13] R. Motwani, A. Naor, and R. Panigrahi, "Lower bounds on locality sensitive hashing," in *Proc. 22nd Annu. ACM Sympos. Comput. Geom.*, 2006, pp. 154–157.
- [14] S. Arya, D. M. Mount, N. S. Netanyahu, R. Silverman, and A. Y. Wu, "An optimal algorithm for approximate nearest neighbor searching in fixed dimensions," *J. Assoc. Comput. Mach.*, vol. 45, no. 6, pp. 891–923, 1998. [Online]. Available: <http://www.cs.umd.edu/mount/Papers/dist.pdf>
- [15] S. Har-Peled, *Geometric Approximation Algorithms*. Amer. Math. Soc., 2011.
- [16] —, "A replacement for Voronoi diagrams of near linear size," in *Proc. 42nd Annu. IEEE Sympos. Found. Comput. Sci.*, 2001, pp. 94–103. [Online]. Available: <http://cs.uiuc.edu/sariel/papers/01/avoronoi/>
- [17] S. Arya and T. Malamatos, "Linear-size approximate Voronoi diagrams," in *Proc. 13th ACM-SIAM Sympos. Discrete Algs.*, 2002, pp. 147–155.
- [18] P. B. Callahan and S. R. Kosaraju, "A decomposition of multidimensional point sets with applications to  $k$ -nearest-neighbors and  $n$ -body potential fields," *J. Assoc. Comput. Mach.*, vol. 42, pp. 67–90, 1995.
- [19] S. Arya, T. Malamatos, and D. M. Mount, "Space-time tradeoffs for approximate nearest neighbor searching," *J. Assoc. Comput. Mach.*, vol. 57, no. 1, pp. 1–54, 2009.
- [20] J. Erickson, "New lower bounds for Hopcroft's problem," *Discrete Comput. Geom.*, vol. 16, pp. 389–418, 1996.
- [21] N. N. Dalvi, C. Ré, and D. Suciu, "Probabilistic databases: Diamonds in the dirt," *Commun. ACM*, vol. 52, no. 7, pp. 86–94, 2009.
- [22] C. C. Aggarwal, *Managing and Mining Uncertain Data*. Springer, 2009.
- [23] P. K. Agarwal, A. Efrat, S. Sankararaman, and W. Zhang, "Nearest-neighbor searching under uncertainty," in *Proc. 31st ACM Sympos. Principles Database Syst.*, 2012, pp. 225–236.
- [24] P. K. Agarwal, B. Aronov, S. Har-Peled, J. M. Phillips, K. Yi, and W. Zhang, "Nearest neighbor searching under uncertainty ii," in *Proc. 32nd ACM Sympos. Principles Database Syst.*, 2013, p. to appear.
- [25] P. K. Agarwal and J. Erickson, "Geometric range searching and its relatives," in *Advances in Discrete and Computational Geometry*, B. Chazelle, J. E. Goodman, and R. Pollack, Eds. Amer. Math. Soc., 1998.
- [26] S. Har-Peled, "Constructing approximate shortest path maps in three dimensions," *SIAM J. Comput.*, vol. 28, no. 4, pp. 1182–1197, 1999.