

Approximating Minimization Diagrams and Generalized Proximity Search*

Sariel Har-Peled[†]

Nirman Kumar[‡]

August 19, 2016

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 multiplicatively weighted points, but the new technique also works for more general distance functions. For example, we get such data-structures for metrics induced by convex bodies, and the nearest furthest-neighbor distance to a set of point sets. Interestingly, our framework also works for distance functions that do not obey the triangle inequality. For many of these functions no near linear size approximation was known before.

1. Introduction

Given a set of functions $\mathcal{F} = \{f_i : \mathbb{R}^d \rightarrow \mathbb{R} \mid i = 1, \dots, n\}$, their minimization diagram 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 $O(n^{1-O(1/d)})$ [Mat92, Cha10]. 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 Sharir and Agarwal [SA95].

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 ,

*Work on this paper was partially supported by NSF AF awards CCF-0915984, CCF-1421231, and CCF-1217462. A preliminary version of this paper appeared in FOCS 2013 [HK13].

[†]Department of Computer Science; University of Illinois; 201 N. Goodwin Avenue; Urbana, IL, 61801, USA; sariel@illinois.edu; <http://sarielhp.org/>.

[‡]Department of Computer Science; University of Illinois; 201 N. Goodwin Avenue; Urbana, IL, 61801, USA; nkumar5@illinois.edu; <http://www.cs.uiuc.edu/~nkumar5/>.

such that given a query point $\mathbf{q} \in \mathcal{X}$, one can find (quickly) the point $\mathbf{n}_{\mathbf{q}} \in \mathbf{P}$ 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.

To see the connection to lower envelopes, consider a set of data points $\mathbf{P} = \{\mathbf{p}_1, \dots, \mathbf{p}_n\}$ in \mathbb{R}^d . Next, consider the set of functions $\mathcal{F} = \{f_1, \dots, f_n\}$, where $f_i(\mathbf{q}) = \|\mathbf{q} - \mathbf{p}_i\|$, 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 $(\mathbf{p}_i, 0)$). Clearly the NN search problem is to evaluate the minimization diagram of these functions at a query point \mathbf{q} , and to find the function f_i which achieves the minimum at \mathbf{q} .

More generally, given a set of n functions, one can think of the functions as defining distances, and the minimization diagram defining a “nearest-neighbor” under those distance functions, by analogy with the above. The nearest-neighbor distance of a query point here, is simply the “height” of the lower envelope at that point.

Exact nearest-neighbor search. The exact nearest-neighbor search 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 [BCKO08]. 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 [Cla88] showed a data-structure with $O(\log n)$ query 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 with the query time [AM93]. Meiser [Mei93] 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 $y \in \mathbf{P}$, such that

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

where $\mathbf{n}_{\mathbf{q}} \in \mathbf{P}$ is the nearest-neighbor to \mathbf{q} in \mathbf{P} . Considerable amount of work was done on this problem, see [Cla06] and references therein.

Indyk and Motwani showed that ANN queries can be reduced to a small number of near-neighbor queries [IM98, HIM12]. 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 $\tilde{O}(n^{1+1/(1+\varepsilon)^2})$ [AI08]. These bounds are near optimal [MNP06].

Low dimensional ANN search. In low dimensions (i.e., \mathbb{R}^d for constant, small d , say $d \leq 20$), one can use linear space (independent of ε) and get ANN query time $O(\log n + 1/\varepsilon^{d-1})$ [AMN⁺98, Har11]. The trade-off for this logarithmic query time is of course an exponential dependence on d . Interestingly, for this data-structure, the approximation parameter ε need not be prespecified during the construction; one can provide it during the query. An alternative approach is to use Approximate Voronoi Diagrams (AVD), introduced by Har-Peled [Har01], 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 [AM02] showed how to build AVD’s of linear complexity (i.e., $O(n/\varepsilon^d)$). Their construction uses Well-Separated Pairs Decomposition [CK95]. Further trade-offs between query time and space for AVD’s were studied by Arya *et al.* [AMM09].

Generalized distance functions: motivation. The algorithms for approximate nearest-neighbor extend to various metrics in \mathbb{R}^d , for example the well known ℓ_p metrics. In particular, previous constructions of AVD’s extend to ℓ_p metrics [Har01, AM02] 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 $f_{\mathbf{p}}(\mathbf{q}) = \omega_{\mathbf{p}} \|\mathbf{q} - \mathbf{p}\|$. 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 . Such a weighted Voronoi diagram is known as the *multiplicative Voronoi diagram*, and even in the plane it can have quadratic complexity. 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 to the client. Of course, this is only one possible distance function, and there are many other such functions that are of interest.

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 computation models) which show that a solution to Hopcroft’s problem requires $\Omega(n^{4/3})$ time [Eri96]. 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.

1.1. Our results

We present conditions that are sufficient to allow efficient approximation of 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 [AM02] 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, but 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 multiplicative 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 $xC + \mathbf{u}$ contains \mathbf{v} , where $xC = \{x\mathbf{p} \mid \mathbf{p} \in C\}$ is the scaling of C by x . So, given a set of n data points $\mathbf{P} = \{\mathbf{p}_1, \dots, \mathbf{p}_n\}$, and n centrally symmetric and fat 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. If the bodies are all fat and of the same size, then intuitively, since the distance functions induced by them, are all “close” to a scaled version of the Euclidean metric, it should be possible to approximate their minimization diagram, as shown for multiplicative weighted Voronoi diagrams. Note however, that we do not assume that the bodies are of similar size, see Section 4.2 for more details.
- (C) **Nearest furthest-neighbor.** Consider a situation where the given input is uncertain; specifically, for the i th *uncertain* point we are given a set of points $\mathbf{P}_i \subseteq \mathbb{R}^d$ where it might lie (the reader might consider the case where the i th point randomly chooses its location out of the points of \mathbf{P}_i). There is a growing interest in how to handle such inputs, as real world measurements are fraught with uncertainty, see [DRS09, Agg09, AESZ12, AAH⁺13] 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 \mathbf{P}_i ; that is, $F_i(\mathbf{q}) = \max_{\mathbf{p} \in \mathbf{P}_i} \|\mathbf{q} - \mathbf{p}\|$. Thus, in the worst case, the nearest uncertain point to the query is at a distance $F(\mathbf{q}) = \min_i F_i(\mathbf{q})$. Using our framework we can approximate this function efficiently, using near linear space, 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. See Section 4.3 for details.

Paper organization. In Section 2 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 2.5 we summarize the results that follow from the framework. In Section 3, we describe the construction of the AVD and its associated data-structure. We describe in Section 4 some specific cases where the new AVD construction can be used. We conclude in Section 5.

2. The framework and summary of results

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.

2.1. Problem statement

We are given a set \mathcal{F} of n functions from \mathbb{R}^d to \mathbb{R}^+ , and a parameter ε . Our task is to build a data-structure of near linear size, such that given a query point \mathbf{q} , one can quickly compute a value x , such that $d(\mathbf{q}) \leq x \leq (1 + \varepsilon)d(\mathbf{q})$, where

$$d(\mathbf{q}) = d(\mathbf{q}, \mathcal{F}) = \min_{i=1, \dots, n} f_i(\mathbf{q}) \quad (2.1)$$

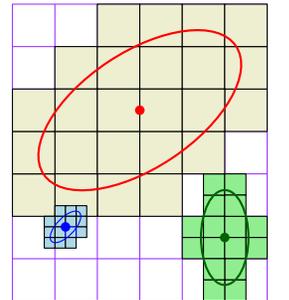
is the lower envelope of \mathcal{F} at \mathbf{q} ; that is, $d(\mathbf{q})$ is the *distance* of \mathbf{q} from \mathcal{F} . Since the functions of \mathcal{F} are distance functions in our applications, we refer to these approximate minimization diagram queries as *approximate nearest neighbor (ANN)* queries.

For such a set of functions \mathcal{F} , an *approximate Voronoi diagram (AVD)* is a decomposition of space into near linear number of cells, each of low descriptive complexity, such that every cell c in the AVD has a small number of functions \mathcal{G}_c associated with it, such that for any point \mathbf{q} , if c is the cell of the AVD that contains \mathbf{q} , then $d(\mathbf{q}, \mathcal{G}_c)$ approximates $d(\mathbf{q}, \mathcal{F})$ well. As such, given an ANN query, it can be answered by doing a point-location query in the AVD to compute c , and then computing $d(\mathbf{q}, \mathcal{G}_c)$ explicitly. Usually, such an AVD is constructed using compressed quadtrees and yield $O(\log n)$ ANN query time, see [Har11] for details.

2.2. Informal description of the technique

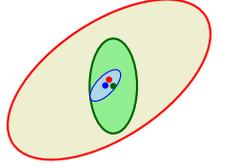
Consider n points in the plane $\mathbf{p}_1, \dots, \mathbf{p}_n$, where the “distance” from the i th point to a query \mathbf{q} , is the minimum scaling of an ellipse \mathcal{E}_i (centered at \mathbf{p}_i), till it covers \mathbf{q} , and let f_i denote this distance function. Assume that these ellipses are fat, i.e., for some constant $\alpha > 0$, there is a ball of some radius r_i , centered at \mathbf{p}_i , that lies inside \mathcal{E}_i , while the ball of radius αr_i covers \mathcal{E}_i . Clearly each function f_i has a graph that is a deformed cone. Given a query point $\mathbf{q} \in \mathbb{R}^2$, we are interested in the first function graph being hit by a vertical ray shot upward from $(\mathbf{q}, 0)$. In particular, let $d(\mathbf{q}) = \min_{i=1, \dots, n} f_i(\mathbf{q})$ be the minimization diagram of these functions.

As a first step in computing $d(\mathbf{q})$, consider the decision version of this problem. Given a value r , we are interested in deciding if $d(\mathbf{q}) \leq r$. That is, we want to decide if $\mathbf{q} \in \bigcup_i (\mathbf{p}_i + r\mathcal{E}_i)$. Of course, this is by itself a computationally expensive task, and as such we satisfy ourselves with an approximate decision procedure. Formally, we replace every ellipse by a collection of grid cells (of the right resolution), such that approximately it is enough to decide if the query point lies inside any of these grid cells – if it does, we know that $d(\mathbf{q}) \leq (1 + \varepsilon)r$, otherwise $d(\mathbf{q}) > r$. Of course, as depicted in the right, since the ellipses are of different sizes, the grid cells generated for each ellipse might belong to different resolutions, and might be of different sizes. Nevertheless, one can perform this point-location query among the marked grid squares quickly using a compressed quadtree.



If we were interested only in the case where $d(\mathbf{q})$ is guaranteed to be in some interval $[\alpha, \beta]$, then the problem would be easily solvable. Indeed, build a sequence of the above deciders $\mathcal{D}_1, \dots, \mathcal{D}_m$, where \mathcal{D}_i is for the distance $(1 + \varepsilon)^i \alpha$, and $m = \log_{1+\varepsilon}(\beta/\alpha)$. Clearly, doing a binary search over these deciders would resolve the distance query.

Sketchable. Unfortunately, in general, there is no such guarantee – this makes the problem significantly more challenging. Fortunately, for truly “large” distances, a collection of such ellipses looks like a constant number of ellipses (at least in the approximate case). In the example of the figure above, for large enough distance, the ellipses looks like a single ellipse, as demonstrated in the figure on the right. Slightly more formally, if $\bigcup_i(\mathbf{p}_i + r\mathcal{E}_i)$ is connected, then the set $\bigcup_i(\mathbf{p}_i + R\mathcal{E}_i)$ can be $(1 + \varepsilon)$ -approximated by a constant number of these ellipses, if $R > \Omega(nr/\varepsilon)$. A family of functions having this property is *sketchable*. This suggests the problem is easy for large distances.



Critical values to search over. The above suggests that connectivity is the underlying property that enables us to simplify and replace a large set of ellipses, by a few ellipses, if we are looking at them from sufficiently far. This implies that the critical values when the level-set of the functions changes its connectivity are the values we should search over during the nearest-neighbor search. Specifically, let r_i be the minimal r when the set $\bigcup_{k=1}^n(\mathbf{p}_k + r_i\mathcal{E}_k)$ has $n - i$ connected components, for $i = 0, 1, \dots, n - 1$, and let $0 = r_0 < r_1 \leq r_2 \leq \dots \leq r_{n-1} < r_n = \infty$, be the resulting sequence. Using the above decision procedure, and a binary search, we can find the largest index j , such that $r_j \leq d(\mathbf{q}) < r_{j+1}$. Furthermore, the decision procedure for the distance r_{j+1} reports which connected components of $\bigcup_{k=1}^n(\mathbf{p}_k + r_j\mathcal{E}_k)$ contain the query point \mathbf{q} in their union, at distance r_{j+1} . Assume the set of these connected components is formed by the first t functions; that is, $\bigcup_{k=1}^t(\mathbf{p}_k + r_{j+1}\mathcal{E}_k)$ contains \mathbf{q} . There are two possibilities:

- (A) If $d(\mathbf{q}) \in [r_j, c_a(t/\varepsilon)r_j]$, then a binary search with the decision procedure would approximate $d(\mathbf{q})$, where c_a is some constant.
- (B) If $d(\mathbf{q}) > (t/\varepsilon)r_j$ then, each connected component of $\bigcup_{k=1}^t(\mathbf{p}_k + r_j\mathcal{E}_k)$ can be sketched and replaced by a constant number of representative functions (the sketch), and the nearest-neighbor search can now be restricted to the union of the sketches. This is an instance of the problem with a smaller number of functions.

2.2.1. Challenges

There are several challenges in realizing the above scheme:

- (A) We are interested in more general distance functions. To this end, we carefully formalize what conditions the underlying distance functions induced by each point have to fulfill, so that our framework applies.
- (B) The above scheme requires (roughly) quadratic space to be realized. To reduce the space to near linear, we need to be more aggressive about replacing clusters of points/functions, by sketches. To this end, we replace our global scheme by a recursive scheme that starts with the “median” critical value, and fork the search at this value using the decision procedure. Now, when continuing the search above this value, we replace every cluster (at this resolution) by its sketch.
- (C) Computing this “median” value directly is too expensive. Instead, we randomly select a function and compute the connectivity radius of this single distance function with the remaining functions. With good probability this value turns out to be good.
- (D) We need to be careful to avoid accumulation in the error as we replace clusters by sketches.

2.3. 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 ℓ , denoted by G_ℓ , is the natural tiling of \mathbb{R}^d , with cubes

of side-length ℓ (i.e., with a vertex at the origin, and axis parallel sides). A cube \square is **canonical** if it belongs to \mathbf{G}_ℓ , ℓ 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 2.1. To approximate a set $X \subseteq [0, 1]^d$, up to distance r , consider the set $\mathbf{G}_{\approx r}(X)$ of all the canonical grid cells of \mathbf{G}_ℓ , that have a non-empty intersection with X , where $\ell = 2^{\lfloor \log_2(r/\sqrt{d}) \rfloor}$. Let $\cup \mathbf{G}_{\approx r}(X) = \bigcup_{\square \in \mathbf{G}_{\approx r}(X)} \square$, denote the union of cubes of $\mathbf{G}_{\approx r}(X)$. Each cube of $\mathbf{G}_{\approx r}(X)$ has diameter at most r .

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

Definition 2.2. For $\ell \geq 0$ and a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, the ℓ **sublevel set** of f is the set $f_{\leq \ell} = \left\{ \mathbf{q} \in \mathbb{R}^d \mid f(\mathbf{q}) \leq \ell \right\}$. The number ℓ will be the **level** of the sublevel set $f_{\leq \ell}$. For a set of functions \mathcal{F} , let $\mathcal{F}_{\leq \ell} = \bigcup_{f \in \mathcal{F}} f_{\leq \ell}$.

For a value r , the approximate versions of the sublevel sets are $f_{\leq \ell, \approx r} = \mathbf{G}_{\approx r}(f_{\leq \ell})$ and $\mathcal{F}_{\leq \ell, \approx r} = \bigcup_{f \in \mathcal{F}} \mathbf{G}_{\approx r}(f_{\leq \ell})$.

Definition 2.3. Given a function f and $\mathbf{q} \in \mathbb{R}^d$, their distance is $\mathfrak{d}(\mathbf{q}, f) = f(\mathbf{q})$. Given two functions f and g , their **distance** $\mathfrak{d}(f, g)$ is the minimum $\ell \geq 0$ such that $f_{\leq \ell} \cap g_{\leq \ell} \neq \emptyset$. Similarly, for two sets of function, \mathcal{F} and \mathcal{G} , their **distance** is $\mathfrak{d}(\mathcal{F}, \mathcal{G}) = \min_{f \in \mathcal{F}, g \in \mathcal{G}} \mathfrak{d}(f, g)$. See also Eq. (2.1).

Example 2.4. To decipher these somewhat cryptic definitions, the reader might want to consider the standard settings of regular Voronoi diagrams. Here, we have a set \mathbf{P} of n points. The i th point $\mathbf{p}_i \in \mathbf{P}$ induces the natural function $f_i(\mathbf{q}) = \|\mathbf{q} - \mathbf{p}_i\|$. We have:

- (A) The graph of f_i in \mathbb{R}^{d+1} is a cone “opening upwards” with an apex at $(\mathbf{p}_i, 0)$.
- (B) The ℓ sublevel set of f_i (i.e., $(f_i)_{\leq \ell}$) is a ball of radius ℓ centered at \mathbf{p}_i .
- (C) The distance of \mathbf{q} from f_i , i.e., $\mathfrak{d}(\mathbf{q}, f_i)$, is the Euclidean distance between \mathbf{q} and \mathbf{p}_i .
- (D) Consider two subsets of points $X, Y \subseteq \mathbf{P}$ and let \mathcal{F}_X and \mathcal{F}_Y be the corresponding sets of functions. The distance $\ell = \mathfrak{d}(\mathcal{F}_X, \mathcal{F}_Y)$ is the minimum radius of balls centered at points of X and Y , such that there are two balls, from the two sets, that intersect; that is, ℓ is half the minimum distance between a point of X and a point of Y . In particular, if the union of balls of radius ℓ centered at X is connected, i.e., $(\mathcal{F}_X)_{\leq \ell}$ is connected, and similarly for Y , then $(\mathcal{F}_X \cup \mathcal{F}_Y)_{\leq \ell}$ is connected. This is the critical value where two connected components of the sublevel sets merge.

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

Observation 2.5. Suppose that f and g are two functions such that $\mathfrak{d}(f, g) > 0$ and $\mathbf{q} \in \mathbb{R}^d$. Then, $\max(\mathfrak{d}(\mathbf{q}, f), \mathfrak{d}(\mathbf{q}, g)) \geq \mathfrak{d}(f, g)$.

Definition 2.6. Let B_1, B_2, \dots, B_m be connected, nonempty sets in \mathbb{R}^d . This collection of sets is **connected** if $\cup_i B_i$ is connected.

Definition 2.7. A family of functions \mathfrak{F} from \mathbb{R}^d to \mathbb{R} , is **well behaved**, if for any m functions $\{f_1, \dots, f_m\} \subseteq \mathfrak{F}$, their lower envelope has the following properties:

- (A) The projection of the lower envelope into \mathbb{R}^d can be decomposed into constant complexity cells. Formally, each such cell is a semialgebraic set that can be described by a constant number of algebraic inequalities, see Basu *et al.* [BPR06].
- (B) One can build a data-structure that answer point-location queries in the resulting set of cells in $O(\log m)$ time.
- (C) The number of cells in this decomposition, the time to compute this decomposition, and the space required to store the data-structure are all bounded by $m^{\gamma \cdot d}$.

The constant γ is the *wellness constant* associated with \mathfrak{F} .

2.3.1. Sketches

A key idea underlying our approach 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 2.8 ($\ell_{\text{con}}(\mathcal{F})$). Given a set of functions \mathcal{G} , if \mathcal{G} contains a single function then the *connectivity level* $\ell_{\text{con}}(\mathcal{G})$ 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 2.6.

Remark. The requirement above that every sublevel set of a function is connected (on its own) in the above definitions can be relaxed. Since we currently we have no application that uses this extra flexibility, we opted to stick with the cleaner and somewhat more restrictive definition – our results should also hold in this more general settings.

Definition 2.10. Given a set of functions \mathcal{G} and $\delta \geq 0, y_0 \geq 0$, a (δ, 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 verify that for any $\mathcal{G}, \delta \geq 0, y_0 \geq 0$, if $\mathcal{H} \subseteq \mathcal{G}$ is a (δ, 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 (δ', y'_0) -sketch for \mathcal{G} . Clearly, the set \mathcal{G} is a sketch of itself.

2.4. Conditions on the functions

A family of functions \mathfrak{F} from \mathbb{R}^d to \mathbb{R}^+ , is *dependable*, if there are absolute constants, $\zeta > 0$ (*growth constant*), and a positive integer c_{sk} (*sketch constant*), depending on \mathfrak{F} , such that the following conditions are satisfied for any $f \in \mathfrak{F}$:

- (P1) **Compactness.** For any $y \geq 0$ the sublevel set $(f)_{\leq y}$ is compact.
- (P2) **Bounded growth.** There is a function $\lambda_f : \mathbb{R}^+ \rightarrow \overline{\mathbb{R}^+}$, called the *growth function*, such that for any $y \geq 0$ and $\varepsilon > 0$, if $f_{\leq y} \neq \emptyset$, then $\lambda_f(y) \geq \text{diam}(f_{\leq y})/\zeta$, 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 $\mathbf{B}(\mathbf{u}, r)$ is the ball of radius r centered at \mathbf{u} .
- (P3) **Existence of a sketch.** Given $\delta > 0$ and a finite subset $\mathcal{G} \subseteq \mathfrak{F}$, there is a $\mathcal{H} \subseteq \mathcal{G}$ with $|\mathcal{H}| = O(1/\delta^{c_{\text{sk}}})$ and $y_0 = O(\ell_{\text{con}}(\mathcal{G}) \cdot (|\mathcal{G}|/\delta)^{c_{\text{sk}}})$ such that, \mathcal{H} is a (δ, y_0) -sketch for \mathcal{G} .

We also require some straightforward properties from the computation model:

- (C1) For all $\mathbf{q} \in \mathbb{R}^d$, the value $f(\mathbf{q}) = d(\mathbf{q}, f)$ is computable in constant time.
- (C2) For any $y \geq 0, r > 0$, the set of grid cells approximating the sublevel set $f_{\leq y}$ of f , that is $f_{\leq y, \approx r} = \mathbf{G}_{\approx r}(f_{\leq y})$ (see Definition 2.1), is computable in linear time in its size.
- (C3) For any $f, g \in \mathfrak{F}$, the distance $d(f, g)$ is computable in constant time.

We also assume that the growth function $\lambda_f(y)$, from Condition (P2), can be computed quickly, i.e., in constant time.

Remark. (A) In the following, the set of functions $\mathcal{F} = \{f_1, \dots, f_n\}$ is taken from an implicit family \mathfrak{F} . Thus, the growth and sketch constants used are associated with \mathfrak{F} and do not depend on n or \mathcal{F} . Henceforth, the terms family, set, or class of functions refers to \mathcal{F} , while the underlying family \mathfrak{F} would be only be implicitly mentioned.

(B) In the following, many of the quantities used above (e.g., resolution r , or various distances used) would be computed approximately, as exact computation is both computationally expensive and not necessary to make things work.

Remark 2.12. Condition (C2) is only used for $r = \Omega(\varepsilon\lambda_f(y))$. That is, the grid used on the sublevel set at a resolution typically ε times its growth function value at its level, which by Condition (C2) is also $\Omega(\varepsilon\text{diam}(f_{\leq y}))$. Here $\varepsilon \in (0, 1)$ is a prespecified approximation parameter. As such, the number of grid cells in the approximation $f_{\leq y, \approx r}$, for a single function f level set, is $O(1/\varepsilon^d)$. The value of r being used can be somewhat imprecise as long as it is sufficiently small.

2.4.1. Properties

In the following, let \mathcal{F} be a family of functions, that satisfies the conditions above. The functions under consideration have the basic properties listed below. Since these properties are straightforward but their proofs are somewhat tedious, and are delegated to Appendix B_{p30}.

- (L1) 0-LEVEL IS EMPTY OR A POINT. For any $f \in \mathcal{F}$, either $f_{\leq 0} = \emptyset$ or $f_{\leq 0}$ consists of a single point, see Lemma B.1_{p30}.
- (L2) CONNECTIVITY LEVEL FOR MORE THAN ONE FUNCTION IS NON-ZERO. If $\ell_{\text{con}}(\mathcal{F}) = 0$ for any non-empty set \mathcal{F} , then $|\mathcal{F}| = 1$. See Definition 2.8 and Observation B.2_{p30}.
- (L3) NEAR CONVEXITY. Let $f \in \mathcal{F}$ and $y \geq 0$. For any points $u, v \in f_{\leq y}$, we have $uv \subseteq f_{\leq (1+\zeta/2)y}$, where uv denotes the segment joining u to v , see Lemma B.3_{p31}.
- (L4) CONNECTIVITY LEVEL IS PRESERVED UNDER SKETCHING. For any $\mathcal{G} \subseteq \mathcal{F}$, $\delta \geq 0$ and $y \geq 0$, such that \mathcal{G} is a (δ, y) -sketch for \mathcal{F} , we have that, $\ell_{\text{con}}(\mathcal{G}) \leq (1 + \delta)(1 + \zeta/2) \max(y, \ell_{\text{con}}(\mathcal{F}))$, see Lemma B.5_{p31}.
- (L5) DISTANCES ARE PRESERVED UNDER SKETCHING. Let $\mathcal{G} \subseteq \mathcal{F}$, such that \mathcal{G} is a (δ, y_0) -sketch for \mathcal{F} , for some $\delta \geq 0$ and $y_0 \geq 0$. Let q be a point such that $d(q, \mathcal{F}) \geq y_0$. Then we have that $d(q, \mathcal{G}) \leq (1 + \delta)d(q, \mathcal{F})$, see Lemma B.6_{p31}.

Remark 2.13 (Computing the sketch). 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 (δ, y_0) -sketch can be computed in time $O(|\mathcal{G}|/\delta^{\text{csk}})$ time. As another example, the minimum of a function $f \in \mathcal{F}$ can be computed efficiently (see (L1)).

2.5. Summary of results

The following is the main result of the paper, see Section 3 for details.

Theorem 2.14. *Let \mathcal{F} be a set of n functions in \mathbb{R}^d that complies with the assumptions of Section 2.4, and has sketch constant $\text{c}_{\text{sk}} \geq d + 1$. Then, one can build a data-structure to answer ANN for this set of functions, with the following properties:*

- (A) *The query time is $O(\log n + 1/\varepsilon^{\text{c}_{\text{sk}}})$.*

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

One can transform the data-structure into an AVD, and in the process improve the query time (the space requirement slightly deteriorates). See Section 3.4.1 for details.

Corollary 2.15. *Consider the settings of Theorem 2.14 with the additional condition that \mathcal{F} is well-behaved, see Definition 2.7, with wellness constant γ . Then, one can build a data-structure to answer ANN for this set of functions, with the following guarantees:*

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

This also computes an AVD of \mathcal{F} of complexity $O(S)$ – a space decomposition into cells, such that each cell has a single function of \mathcal{F} associated with it, and for any point in this cell, this function is a $(1 + \varepsilon)$ -ANN among the functions of \mathcal{F} .

2.5.1. Distance functions for which the framework applies

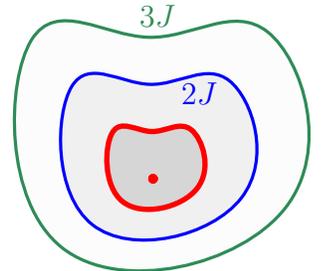
2.5.1.1. Multiplicative distance functions with additive offsets. For $i = 1, \dots, n$, let \mathbf{p}_i be a point with an associated weight $w_i > 0$, and an **offset** $\alpha_i \geq 0$. The **multiplicative distance with offset** induced by the i th point is $f_i(\mathbf{q}) = w_i \|\mathbf{q} - \mathbf{p}_i\| + \alpha_i$. Section 4.1 shows that these distance functions satisfy the conditions of Section 2.4, implying the following result.

Theorem 2.16. *Consider a set \mathbf{P} of n points in \mathbb{R}^d , where the i th point \mathbf{p}_i has additive weight $\alpha_i \geq 0$ and multiplicative weight $w_i > 0$, for $i = 1, \dots, n$. Such a point induces the additive/multiplicative distance function $f_i(\mathbf{q}) = w_i \|\mathbf{q} - \mathbf{p}_i\| + \alpha_i$. One can compute a $(1 + \varepsilon)$ -AVD for these distance functions, with near linear space complexity, and logarithmic query time. See Theorem 2.14_{p9} and Corollary 2.15 for the precise bounds, with the sketch and wellness constants $c_{\text{sk}} = d + 1$ and $\gamma = 1$.*

2.5.1.2. Scaling distance. Somewhat imprecisely, a connected body J centered at a point ρ is α -rounded fat if it is α -fat (that is, there is radius r such that $\text{ball}(\rho, r) \subseteq J \subseteq \text{ball}(\rho, \alpha r)$), and from any point \mathbf{p} on the boundary of J the “cone” $\mathcal{CH}(\text{ball}(\rho, r) \cup \mathbf{p})$ is contained inside J (i.e., every boundary point sees a large fraction of the “center” of the object). Furthermore, assume that the boundary of each object J has constant descriptive complexity. Note that such an object is star-shaped with ρ as its center. See Definition 4.5_{p22} for the formal definition.

For such an object J , its **scaling distance** to a point \mathbf{q} , is the minimum t , such that $\mathbf{q} \in tJ$ (where the scaling of J is done around its center ρ), see figure on the right. Given n α -rounded fat objects, it is natural to ask for the Voronoi diagram induced by their scaling distance.

The distance functions induced by such a set of objects complies with the framework of Section 2.4 implying the following result, see Section 4.2 for details.



Theorem 2.17. *Consider a set \mathcal{J} of n α -rounded fat objects in \mathbb{R}^d , for some constant α . Then one can compute a ANN data-structure, for the scaling distance functions induced by \mathcal{J} , with near linear space complexity, and logarithmic query time. See Theorem 2.14_{p9} for the precise bounds, with the sketch constant being $c_{\text{sk}} = d + 1$.*

2.5.1.3. 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 $f_S(\mathbf{q}) = \max_{\mathbf{s} \in S} \|\mathbf{q} - \mathbf{s}\|$; 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 $f_S(\cdot)$ induces a partition of space into 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 sized representation (in constant dimensions), see [Har99].

Given n sets of points P_1, \dots, P_n in \mathbb{R}^d , we are interested in the distance function $f(\mathbf{q}) = \min_i f_i(\mathbf{q})$, where $f_i(\mathbf{q}) = f_{P_i}(\mathbf{q})$. This quantity arises naturally when one tries to model uncertainty [AAH⁺13]; 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, $f_i(\mathbf{q})$ is the worst case distance to the i th point, and $f(\mathbf{q})$ is the worst-case nearest-neighbor distance to the random point-set generated by picking the i th point from P_i , for $i = 1, \dots, n$. The function $f(\cdot)$ is the *nearest furthest-neighbor* distance, and the task at hand is to approximate it.

The distance functions f_1, \dots, f_n satisfy the conditions of the framework, see Section 4.3, implying the following.

Theorem 2.18. *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 these points into a data-structure for answering $(1 + \varepsilon)$ -ANFN (approximate nearest further-neighbor) queries, with the following guarantees:*

(A) *The query time is $O(n + 1/\varepsilon^{d+1})$.*

(B) *The preprocessing time and space used is $O\left(n(\varepsilon^{-1} \log n)^{O(d)}\right)$.*

Note that the space and query time used, depend only on n , and not on the input size.

3. Constructing the AVD

The input is a set \mathcal{F} of n functions, satisfying the conditions of Section 2.4, and a number $\varepsilon \in (0, 1)$. We preprocess \mathcal{F} , such that given a query point \mathbf{q} , one can compute a $f \in \mathcal{F}$, where $d(\mathbf{q}, \mathcal{F}) \leq d(\mathbf{q}, f) \leq (1 + \varepsilon)d(\mathbf{q}, \mathcal{F})$.

3.1. Building blocks

3.1.1. 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 $\mathcal{D}_{near} = \mathcal{D}_{nr}(\mathcal{G}, \varepsilon, \alpha)$ can decide (approximately) if a point has distance larger or smaller than α from \mathcal{G} . 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 $d(\mathbf{q}, \mathcal{G}) \in (\alpha, (1 + \varepsilon)\alpha]$. 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 $m = |\mathcal{G}|$.

Lemma 3.1. *Given a set $\mathcal{G} \subseteq \mathcal{F}$ of m functions, $\alpha > 0$ 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 α , in time $T_{\leq}(m) = O(\log(m/\varepsilon))$.*

Proof: For each $f \in \mathcal{G}$, consider the canonical grid set $\mathbf{G}_{\approx r_f}(f_{\leq \alpha})$, where $r_f = \varepsilon \lambda_f(\alpha) \geq \varepsilon \text{diam}(f_{\leq \alpha})/\zeta$, where $\lambda_f(\cdot)$ and ζ , are the growth function and the growth constant, respectively, see (P2)_{p8}. The sublevel set of interest is $\mathcal{G}_{\leq \alpha}$, and its approximation is $\mathcal{C} = \bigcup_{f \in \mathcal{G}} \mathbf{G}_{\approx r_f}(f_{\leq \alpha})$, as the bounded growth condition (P2)_{p8} implies that $f_{\leq \alpha} \subseteq \mathbf{G}_{\approx r_f}(f_{\leq \alpha}) \subseteq f_{\leq (1+\varepsilon)\alpha}$. The set of canonical cubes \mathcal{C} can be stored in a compressed quadtree \mathcal{T} , and given a query point, we can decide if a point is covered by some cube of \mathcal{C} , by performing a point location query in \mathcal{T} .

By Remark 2.12, $|\mathbf{G}_{\approx r_f}(f_{\leq \alpha})| = O(\varepsilon^{-d})$. As such, the total number of canonical cubes in \mathcal{C} is $O(m/\varepsilon^d)$, and the compressed quadtree for storing them can be computed in $O(m\varepsilon^{-d} \log(m/\varepsilon))$ time [Har11].

We mark a cell of the resulting quadtree by the function whose sublevel set it arose from (ties can be resolved arbitrarily). During query, if \mathbf{q} is found in one of the cells we return **yes** and the function associated with the cell, otherwise we return **no**.

If we have that $\mathfrak{d}(\mathbf{q}, \mathcal{G}) \leq \alpha$, then the query point \mathbf{q} will be found in one of the marked cells, since they cover $\mathcal{G}_{\leq \alpha}$. As such, the query will return **yes**. Moreover, if the query does return a **yes**, then it belongs to a cube of \mathcal{C} that is completely covered by $\mathcal{G}_{\leq (1+\varepsilon)\alpha}$, as desired. \blacksquare

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

The time to perform an interval query is denoted by $T_I(m, \alpha, \beta)$.

Lemma 3.2. *Given a set \mathcal{G} of m functions, an interval $[\alpha, \beta]$, and an approximation parameter $0 < \tau \leq 4$, such that $\log(4\beta/\alpha) = O(m/\tau)$, one can construct an interval data-structure of size $O(m\tau^{-d-1} \log \frac{4\beta}{\alpha})$, in time $O(m\tau^{-d-1} \log \frac{4\beta}{\alpha} \log \frac{m}{\tau})$, such that given a query point \mathbf{q} one can answer a $(1 + \tau)$ -ANN query for the distances in the interval $[\alpha, \beta]$, in time $T_I(m, \alpha, \beta) = O\left(\log \frac{m \log(4\beta/\alpha)}{\tau}\right)$.*

Proof: Using Lemma 3.1, build a $(1 + \tau/4)$ near-neighbor data-structure \mathbf{D}_i for \mathcal{G} , for distance $r_i = (\alpha/2)(1 + \tau/4)^i$, for $i = 0, \dots, L = \lceil \log_{1+\tau/4}(4\beta/\alpha) \rceil = O(\tau^{-1} \log(4\beta/\alpha))$. Clearly, an interval query can be answered in three stages:

- (A) Perform a point-location query in \mathbf{D}_0 . If the answer is **yes** then $\mathfrak{d}(\mathbf{q}, \mathcal{G}) < \alpha$. We can also return a function $g \in \mathcal{G}$ with $\mathfrak{d}(\mathbf{q}, g) < \alpha$.
- (B) Similarly, perform a point-location query in \mathbf{D}_L . If the answer is **no** then $\mathfrak{d}(\mathbf{q}, \mathcal{G}) > \beta$ and we are done.
- (C) It must be that $\mathfrak{d}(\mathbf{q}, \mathcal{G}) \in [r_i, r_{i+1}]$ for some i . Find this i by performing a binary search on the data-structures $\mathbf{D}_0, \dots, \mathbf{D}_L$, for the first i such that \mathbf{D}_i returns **no**, but \mathbf{D}_{i+1} returns **yes**. Clearly, \mathbf{D}_{i+1} provides us with the desired $(1 + \tau/4)^2$ -ANN to the query point.

To get the improved query time, observe that we can overlay these compressed quadtrees $\mathbf{D}_0, \dots, \mathbf{D}_L$, into a single quadtree. For every leaf (or compressed node) of this quadtree, we compute the original node with the lowest value covering this node. Clearly, finding the desired distance can now be resolved by a single point-location query in this overlay of quadtrees. The total size of these quadtrees is

$S = O(Lm/\tau^d)$, and the total time to compute these quadtrees is $T_1 = O(L(m/\tau^d) \log(m/\tau))$, and the time to compute their overlay is

$$O(S \log L) = O\left(\frac{Lm}{\tau^d} \log L\right) = O\left(\frac{m\tau^{-1} \log(4\beta/\alpha)}{\tau^d} \log(m\tau^{-1} \log(4\beta/\alpha))\right) = O\left(\frac{m \log(4\beta/\alpha)}{\tau^{d+1}} \log \frac{m}{\tau}\right),$$

since $\log(4\beta/\alpha) = O(m/\tau)$. Indeed, overlaying L quadtrees is equivalent to merging L sorted lists, see [Har11, Lemma 2.19]. The time to perform a point-location query in the overlaid quadtree is $O(\log S)$. ■

Lemma 3.2 readily implies that if the nearest-neighbor distance lies in an interval of values of polynomial spread, then Lemma 3.2 yields the desired data-structure. To overcome this unbounded spread problem, we first argue that under our assumptions there are only linear number of intervals where interesting things happen to the distance function.

3.1.3. Connected components of the sublevel sets

Given a finite set X and a partition 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 $i = 1, \dots, k$, each X_i is a **part** of the partition.

Definition 3.3. 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 .

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 Π whose union is $X'_i \in \Xi$.

Observation 3.4. (A) Given partitions Π and Ξ of a finite set X , if $\Pi \sqsubseteq \Xi$ then $|\Xi| \leq |\Pi|$.

(B) 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 i , the set $\phi(\Pi, \Xi, i)$ can be returned in $O(|\phi(\Pi, \Xi, i)|)$ time, and its size can be computed in constant time.

Definition 3.5. Given partitions $\Pi = \langle X_1, \dots, X_k \rangle_X \sqsubseteq \Xi = \langle X'_1, \dots, X'_{k'} \rangle_X$ of a set X , for any part X'_i of Ξ , let $\Pi[X'_i]$ denote the **induced partition** of X'_i by Π , where $\Pi[X'_i] = \langle X_{i_1}, \dots, X_{i_r} \rangle_{X'_i}$, and $\{i_1, \dots, i_r\} = \phi(\Pi, \Xi, i)$.

Definition 3.6. 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 ℓ . And the partition of \mathcal{G} by these clusters, denoted by $\mathcal{C}(\mathcal{G}, \ell)$, is the **ℓ -clustering** of \mathcal{G} .

The values ℓ at which the ℓ -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 the key in decomposing the nearest-neighbor search on \mathcal{F} , into a search on smaller sets.

Observation 3.7. If $0 \leq a \leq b$, then $\mathcal{C}(\mathcal{G}, a) \sqsubseteq \mathcal{C}(\mathcal{G}, b)$.

For any number $\ell \geq 0$, the following lemma testifies that ℓ -clustering can be approximated quickly.

Lemma 3.8. 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 $\mathcal{C}(\mathcal{G}, \ell) \sqsubseteq \Psi \sqsubseteq \mathcal{C}(\mathcal{G}, (1 + \varepsilon)\ell)$, where $m = |\mathcal{G}|$.

Proof: For each $f \in \mathcal{G}$, tile the sublevel sets $(f)_{\leq \ell}$ by canonical cubes of small enough diameter, such that the bounded growth condition (P2)_{p8} assures that all the cubes are inside $(f)_{\leq (1+\varepsilon)\ell}$. To this end, for $f \in \mathcal{G}$, set $r_f = \varepsilon \lambda_f(\ell) \geq \varepsilon \text{diam}(f_{\leq \ell})/\zeta$, and compute the set $\mathcal{C}_f = f_{\leq \ell, \approx r_f}$, see Definition 2.2_{p7}. It is easy to verify that

$$(\mathcal{G})_{\leq \ell} \subseteq \bigcup_{f \in \mathcal{G}} (\cup \mathcal{C}_f) \subseteq (\mathcal{G})_{\leq (1+\varepsilon)\ell}, \quad (3.1)$$

see Definition 2.1. By assumption, we have that $|\mathcal{C}_f| = O(1/\varepsilon^d)$, and the total number of canonical cubes, in all the sets \mathcal{C}_f for $f \in \mathcal{G}$, is $O(m/\varepsilon^d)$. We throw all these canonical cubes into a compressed quadtree, this takes $O((m/\varepsilon^d) \log(m/\varepsilon))$ time. Here, every node of the compressed quadtree is marked if it belongs to some of these sets, and if so, to which of the sets. Two sets $\cup \mathcal{C}_f$ and $\cup \mathcal{C}_g$ intersect, if and only if there are two canonical cubes, in these two sets, such that they overlap; that is, one of them is a sub-cube of the other. Initialize a union-find data-structure, and traverse the compressed quadtree using **DFS**, keeping track of the current connected component, and performing a union operation whenever encountering a marked node (i.e., all the canonical nodes associated with it, are unionized into the current connected component). Finally, we perform a union operation for all the cells in \mathcal{C}_f , for all $f \in \mathcal{G}$. Clearly, this results in the desired connected components of the intersection graph of $\cup \mathcal{C}_f$ (note that we consider two sets as intersecting only if their interiors intersect). Translating each such connected set of canonical cubes back to the functions that gave rise to them, results in the desired partition. \blacksquare

The partition Ψ computed by Lemma 3.8 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 $\ell_{\text{con}}(C) \leq (1 + \varepsilon)\ell$.

3.1.4. Computing a splitting distance

Definition 3.9. Given a partition $\Psi = \Psi_1(\mathcal{G}, \ell)$ of \mathcal{G} , with $m = |\Psi|$ clusters, a distance x is a *splitting distance* if $m/4 \leq |\Psi_1(\mathcal{G}, x/4)|$, $\Psi \sqsubseteq \Psi_1(\mathcal{G}, x/4)$, and $|\Psi_1(\mathcal{G}, x)| \leq (7/8)m$.

Lemma 3.10. *Given a partition $\Psi = \Psi_1(\mathcal{G}, \ell)$ 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 Ψ .*

Proof: For each cluster $C \in \Psi$, let r_C be its distance from all the functions in $\mathcal{G} \setminus C$; that is $r_C = \min_{f \in C} \min_{g \in \mathcal{G} \setminus C} d(f, g)$. By Lemma 3.8, any two functions that are in distance $\leq \ell$ are in the same cluster of Ψ . This implies that $r_C > \ell$. Now, let $r_1 \leq r_2 \leq \dots \leq r_m$, be these distances for the m clusters of Ψ . We randomly pick a cluster $C \in \Psi$ and compute r_C for it, by brute force – computing the distance of each function of C with the functions of $\mathcal{G} \setminus C$.

Let $\ell' = \max(r_C, 4\ell)$. Let i be the rank of 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, to the nearest cluster other than itself, is larger than $\ell'/4$. We conclude that $|\Psi_1(\mathcal{G}, \ell'/4)| \geq m/4$.
- (B) All the clusters of Ψ 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 $m' \geq m/4$ 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 $m - m'$,

```

AprxNN(  $\mathcal{G}, \Upsilon, \mathbf{q}$  )
  //  $\mathcal{G}$ : set of functions,  $\Upsilon = \Psi_1(\mathcal{G}, \ell)$  for some value  $\ell$ ,  $\mathbf{q}$ : query point.
  if  $|\Upsilon| \leq \log n$  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 3.10.

  // Perform an interval ANN query on  $[x/8, Nx]$ , see Lemma 3.2.
  if  $d(\mathbf{q}, \mathcal{G}) \in [x/8, Nx]$  or  $(1 + \varepsilon/4)$ -ANN found then
    return nearest function found by  $(1 + \varepsilon/4)$ -ANN interval query.
  if  $d(\mathbf{q}, \mathcal{G}) < x/8$  then
     $f \leftarrow$  2-approximate NN query on  $\mathcal{G}$ , with distance  $x/8$ , see Lemma 3.1.
    Find cluster  $C \in \Psi_1(\mathcal{G}, x/4)$ , such that  $f \in C$ , see Lemma 3.8.
    return AprxNN(  $C, \Upsilon[C], \mathbf{q}$  )

  // Must be that  $d(\mathbf{q}, \mathcal{G}) > Nx$ 
  return AprxNN(  $\mathcal{G}, \Psi_1(\mathcal{G}, x), \mathbf{q}$  ) (★)

```

Figure 3.1: Search algorithm: 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 Υ being the partition of \mathcal{F} into singletons (i.e., $\ell = 0$).

while the number of components of size at least 2, is at most $m'/2$. It follows that the total number of components is at most $m - m'/2 \leq 7m/8$. Since each such component corresponds to a cluster in $\Psi_1(\mathcal{G}, \ell')$, the claim is proved.

Now, compute $\Psi_1(\mathcal{G}, \ell')$ and $\Psi_1(\mathcal{G}, \ell'/4)$, using Lemma 3.8. It is clear that, $\Psi \sqsubseteq \Psi_1(\mathcal{G}, \ell'/4)$. 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. \blacksquare

3.2. The search procedure

3.2.1. An initial “naive” implementation

The search procedure is presented in Figure 3.1.

Lemma 3.11. **AprxNN**($\mathcal{G}, \Upsilon, \mathbf{q}$) 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 **AprxNN** is $h = O(\log n)$, where $n = |\mathcal{G}|$.

Proof: The proof is by induction on the size of Υ . If $|\Upsilon| \leq \log n$, then the function realizing $d(\mathbf{q}, \mathcal{G})$ is returned, and the claim is true.

Let x be the computed splitting distance of Υ . The procedure performs an $(1 + \varepsilon/4)$ -approximate interval nearest-neighbor query for \mathbf{q} on the range $[x/8, Nx]$. If this computed the approximate nearest-neighbor, then we are done.

Otherwise, it must be that either $d(\mathbf{q}, \mathcal{G}) < x/8$ or $d(\mathbf{q}, \mathcal{G}) > Nx$, and significantly, we know which of the two options it is:

- (A) If $d(\mathbf{q}, \mathcal{G}) < x/8$, then doing an approximate near-neighbor query on \mathcal{G} , and distance $x/8$, returns a function $f \in \mathcal{G}$ such that $d(\mathbf{q}, f) < x/4$. Clearly, the nearest-neighbor to \mathbf{q} must be in the

cluster containing f in the partition $\Psi_1(\mathcal{G}, x/4)$, and **AprxNN** recurses on this cluster. Now, by induction, the returned ANN is correct.

Since x is a splitting distance for Υ , see Definition 3.9, we have $|\Upsilon|/4 \leq |\Psi_1(\mathcal{G}, x/4)|$ and $\Upsilon \sqsubseteq \Psi_1(\mathcal{G}, x/4)$. As such, since C is one of the clusters of $\Psi_1(\mathcal{G}, x/4)$, the induced partition of C by Υ (i.e., $\Upsilon[C]$), can have at most $(1 - 1/4)|\Upsilon| + 1 \leq (7/8)|\Upsilon|$ clusters.

- (B) Otherwise, we have $d(\mathbf{q}, \mathcal{G}) > Nx$. Since x is a splitting distance, we have that $|\Psi_1(\mathcal{G}, x)| \leq (7/8)|\Upsilon|$, see Definition 3.9. We recurse on \mathcal{G} , and a partition that has fewer clusters, and by induction the returned answer is correct.

In each step of the recursion, the number of sets in the partition shrunk by at least a fraction of $7/8$. As such, after a logarithmic number of recursive calls, the procedure is done. \blacksquare

Remark 3.12. Case (B) in the proof of Lemma 3.11 does not introduce any error, because we are still conceptually working with all the functions. For a more efficient implementation, we would have to sketch the sets of functions, and that would introduce an error in the search which can accumulate. See Corollary 3.13 below.

3.2.2. But where is the beef? Modifying **AprxNN** to provide fast query time

The reader might wonder how we are going to get an efficient algorithm out of **AprxNN**, as in the case that Υ has a small number of clusters (i.e., $\log n$), still requires us to perform a scan on all the functions in these clusters, and compute their distance from the query point \mathbf{q} . It turns out, that we can assume that each cluster of Υ has size bounded by $O(1/\varepsilon^{c_{\text{sk}}})$, i.e., it is really a sketch (of some set of functions), and all the functions in the final set of at most $\log n$ clusters, have also been sketched to form a sketch of size $O(1/\varepsilon^{c_{\text{sk}}})$. Initially, since all the clusters are singletons, this is clearly true. The merge of clusters merge happens during step (\star) of **AprxNN**, see Figure 3.1. In future iterations one can pass the sketches as proxies for the actual cluster, bounding the size of each cluster. In particular, the query time is $O(1/\varepsilon^{c_{\text{sk}}} + \log^2 n)$. 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.

3.2.3. Sketching a sketch

To improve the efficiency of the preprocessing for **AprxNN**, we are going to use sketches more aggressively. Specifically, for each of the clusters of Υ , we can compute their δ -sketches, for $\delta = \varepsilon/(8h) = O(\varepsilon/\log n)$, see Lemma 3.11. From this point on, when we manipulate this cluster, we do it on its sketch. To make this work set $N = n^{4c_{\text{sk}}}$. Indeed, for any set of functions \mathcal{G} under consideration, their δ -sketch is active at a distance $O(\ell_{\text{con}}(\mathcal{G})(|\mathcal{G}|/\delta)^{c_{\text{sk}}}) = O(\ell_{\text{con}}(\mathcal{G})(n/\delta)^{c_{\text{sk}}})$, see (P3)_{p8}. Now assuming $\delta = O(\varepsilon/\log n)$, and ε is at least $1/\log n$ (otherwise our algorithm is no worse than brute-force search, essentially), as well as n sufficiently large, it follows that a sketch is active at distance $\geq N\ell_{\text{con}}(\mathcal{G})$.

The only place in the algorithm where we need to compute the sketches, is in step (\star) of **AprxNN**, see Figure 3.1. Specifically, we compute $\Psi_1(\mathcal{G}, x)$, and for each new cluster $C \in \Psi_1(\mathcal{G}, x)$, 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 δ -sketch for this set, and this serves as this cluster's sketch from this point on. In particular, the recursive calls to **AprxNN** 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 large enough to be 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 connectivity level of the clusters being merged. Indeed, observe that $d(\mathbf{q}, \mathcal{G}) > Nx$ and Nx is N times bigger than x , which is an upper bound on the connectivity level of the clusters being merged.

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 $(1 + \delta)^h \leq 1 + \varepsilon/4$.

If the number of clusters is small enough, the recursion bottoms out; here we store a $\varepsilon/8$ -sketch of the set of functions in the union of the parts, in the search structure. This would reduce the sketch size to $O(1/\varepsilon^{c_{sk}})$. Note however, that this still does not help us as far as recursion - we must pass the larger δ -sketches in the recursive call in step (\star) of **AprxNN**, see Figure 3.1.

Corollary 3.13. *The modified search procedure that uses sketches, as described above, returns a $(1 + \varepsilon)$ -ANN to the query point.*

Proof: The accumulated error in the NN search because of the sketching is bounded by $1 + \varepsilon/4$, by the above. Another factor of $1 + \varepsilon/4$ is introduced by the interval data-structure used by **AprxNN**, if the interval data-structure finds the desired nearest neighbor. Finally, the shrinking of the sketch at a leaf node introduces another $1 + \varepsilon/8$ factor to the error.

Now, following the proof of Lemma 3.11 implies that the error is bounded by $(1 + \varepsilon/4)(1 + \varepsilon/4)(1 + \varepsilon/8) \leq 1 + \varepsilon$. ■

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.

3.3. The connectivity tree, and the preprocessing

Given a set of functions \mathcal{F} , consider the tree tracking the connected components of the sublevel sets of the functions, as the level changes continuously. 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, during the execution of Kruskal's algorithm for MST, for the graph on the functions where distances are defined using $d(\cdot, \cdot)$.

To make the presentation consistent, let $d_{\approx}(X, Y)$ be the minimum x , such that $\Psi_1(X \cup Y, x)$ is connected (see Lemma 3.8_{p13}), and x is a power of two. 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 x is between $\alpha/2$ and α . In particular, since x must be a power of two, there are at most 2 candidate values to consider, and which is the right one can be decided using Lemma 3.8.

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 ℓ , let $L_{\mathcal{H}}(\ell)$ be the set of nodes such that their label is smaller than ℓ , but their parent label is larger than ℓ . 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. Recall that we are using the approximate distances $d_{\approx}(X, Y)$ between sets, and all splitting distances, as well as $N = n^{4c_{\text{sk}}}$ are set to be the closest power of 2.

Lemma 3.14. *Consider a recursive call $\text{AprxNN}(\mathcal{G}, \Upsilon, \mathbf{q})$, see Figure 3.1_{p15}, made during the search algorithm execution. Then there exists a node u , and a value ℓ , such that $\mathcal{G} = F(u)$ 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 AprxNN corresponds to a subtree of \mathcal{H} .*

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

Specifically, the set of all functions stored in a tree T is the set $F(T) = \bigcup_{u \in T} F(u)$, which is equal to $F(\text{root}(T))$. On the other hand, this set of functions is also $F(T) = \bigcup_{u \text{ is a leaf of } T} F(u)$ (which is the most refined partition associated with this tree). Namely, a tree (informally) encodes different partitions that are all refinements of $F(T)$, that (somewhat imprecisely) corresponds to different resolutions. Our task is to find the right resolution/partition where we can readily answer the ANN query, and this is done by doing a divide-and-conquer on the tree T . In the process, the set of functions under consideration would shrink (i.e., as the root of the current tree changes to a subtree's root), and the size of the set of possible partitions would become smaller, as the number of leafs in the subtree similarly decreases. As such, either we would find the ANN during this process, or we would be left with a single node, where the ANN query can be answered directly using the sketch.

Reinterpreting AprxNN (see Figure 3.1) as working on subtrees, results in the following algorithm:

- (A) If T is a single node u , then find the closest function in $F(u)$ to \mathbf{q} . 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 leafs of T .
Note that the number of clusters associated with the finest partition induced by T is exactly the number of leafs of T . As such, the splitting distance (see Lemma 3.10), is exactly such a resolution that shrinks the numbers of leafs in a subtree, by a constant factor.
- (C) Using interval data-structure determine if the distance $d(\mathbf{q}, F(T))$ is in the range $[x/8, Nx]$. If so, we found the desired ANN.
- (D) If $d(\mathbf{q}, F(T)) > Nx$ then continue recursively on portion of T above $L_T(x)$.
- (E) If $d(\mathbf{q}, F(T)) < x/8$ then we know the node $u \in L_T(x/4)$ such that the ANN belongs to $F(u)$. Continue the search recursively on the subtree of T rooted at u .

That is, AprxNN breaks T into subtrees, and continues the search recursively on one of the subtrees. Significantly, every such subtree has size which is at most a 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 AprxNN . 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, and we only use the above \mathcal{H} in analyzing the preprocessing running time.

In particular, given a subtree T with m edges, the corresponding partition Υ would have at most m sets. Each such set would have a δ -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 takes $U_1 = O(M/\varepsilon^{c_{\text{sk}}}) = O(m(\varepsilon\delta)^{-c_{\text{sk}}})$ time, see Remark 2.13_{p9}. Computing the splitting distance, using Lemma 3.10, takes $U_2 = O(M(\log M + 1/\delta^{c_{\text{sk}}})) = O(m \log n \delta^{-c_{\text{sk}}} + m \delta^{-2c_{\text{sk}}})$

time. Computing the interval data-structure of Lemma 3.2 takes

$$U_3 = O(M\varepsilon^{-d-1} \log n \log M) = O(m\delta^{-c_{\text{sk}}}\varepsilon^{-d-1} \log^2 n)$$

time, and requires $S_1 = O(M\varepsilon^{-d-1} \log n) = O(m\delta^{-c_{\text{sk}}}\varepsilon^{-d-1} \log n)$ space. here, the required condition on β, α and m in Lemma 3.2 holds as **AprxNN** bottoms out when the set of functions is of size $O(\log n)$. 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 than the original tree). Finally, we need to compute the δ -sketches for the clusters sent to the appropriate recursive calls, and this takes $U_4 = O(M/\delta^{c_{\text{sk}}}) = O(m\delta^{-2c_{\text{sk}}})$ time, by Remark 2.13_{p9}.

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

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

Now, it is valid to assume a larger value for the sketch constant c_{sk} , than its optimal value as per its definition by (P3)_{p8}. As such we assume that $c_{\text{sk}} \geq (d+1) \geq 2$, and we can simplify the above expression to $O(\varepsilon^{-2c_{\text{sk}}} \log^{2c_{\text{sk}}} n)$. 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 argument, that each edge requires $O(\log n \cdot (S_1/m)) = O(\varepsilon^{-d-1-c_{\text{sk}}} \log^{c_{\text{sk}}+1} n)$. Moreover, the space used to store the $\varepsilon/8$ -sketches at the leaf nodes is $O(n\varepsilon^{-c_{\text{sk}}})$. As such, the overall space used by the data-structure is $(n\varepsilon^{-d-1-c_{\text{sk}}} \log^{c_{\text{sk}}+1} n)$. 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.

3.4. The result

Proof of Theorem 2.14_{p9}: The query time stated above is $O(\log^2 n + 1/\varepsilon^{c_{\text{sk}}})$. To get the improved query time, we observe that **AprxNN** 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. Such an ordering can be assigned, as, conceptually, the search proceeds using queries on the nodes of a tree, down from the root. 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). Furthermore, since not all query points can be dealt with by one of the interval near-neighbor data-structures, we must store sketches of size $O(1/\varepsilon^{c_{\text{sk}}})$, for certain regions of $[0, 1]^d$ – such regions are not covered by any of the cubes arising from the interval near-neighbor data-structures.

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

3.4.1. The AVD construction

Proof of Corollary 2.15_{p10}: We build the data-structure of Theorem 2.14, except that instead of linearly scanning the sketch at a leaf node, during the query time, we preprocess each such sketch (which is set of functions of size $O(1/\varepsilon^{c_{\text{sk}}})$), for an exact point-location query; that is, we compute the lower envelope of the sketch and preprocess it for vertical ray shooting. This can be done because we are assuming that our function family satisfies the well-behavedness condition, see Definition 2.7. This would require $O(1/\varepsilon^{\gamma d c_{\text{sk}}})$ space and time, and the linear scanning that takes $O(1/\varepsilon^{c_{\text{sk}}})$ time, now is replaced by a point-location query that takes $O(\log 1/\varepsilon)$, as desired. The total query time is thus $O(\log n + \log 1/\varepsilon) = O(\log n/\varepsilon)$.

Using a similar argument, the space requirement is $O(n\varepsilon^{-d-1-c_{\text{sk}}} \log^{c_{\text{sk}}+1} n + n\varepsilon^{-\gamma d c_{\text{sk}}})$, where the first part comes from Theorem 2.14, and for the second part we notice that we replaced the space used for sketches, $O(n\varepsilon^{-c_{\text{sk}}})$ by $O(n\varepsilon^{-\gamma d c_{\text{sk}}})$, the space required to store the point-location data structures defined over the decomposition into cells of the projection of the lower envelope of the sketch onto \mathbb{R}^d . One need to clip each such cell to the region associated with the node/leaf (of the compressed quadtree) that uses it, but since each such cell has constant descriptive complexity as a semi-algebraic set, this clipping does not change its asymptotic complexity. \blacksquare

4. Applications

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

4.1. 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 4.2. Here the sublevel sets look like expanding 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 4.2 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 $P = \{\mathbf{p}_1, \dots, \mathbf{p}_n\}$. For $i = 1, \dots, n$, the point \mathbf{p}_i has weight $w_i > 0$, and a constant $\alpha_i \geq 0$, associated with it. We define $f_i(\mathbf{q}) = w_i \|\mathbf{q} - \mathbf{p}_i\| + \alpha_i$. Let $\mathcal{F} = \{f_1, \dots, f_n\}$. We have, $(f_i)_{\leq y} = \emptyset$ for $y < \alpha_i$, and $(f_i)_{\leq y} = B\left(\mathbf{p}_i, \frac{y - \alpha_i}{w_i}\right)$ for $y \geq \alpha_i$. Checking conditions (C1)_{p8} and (C2)_{p8} is trivial. As for (C3)_{p8} we have the following easy lemma,

Lemma 4.1. *For any i, j , such that $1 \leq i, j \leq n$, we have $d(f_i, f_j) = \max\left(\alpha_i, \alpha_j, \|\mathbf{p}_i - \mathbf{p}_j\| \frac{w_i w_j}{w_i + w_j} + \frac{\alpha_i w_j + \alpha_j w_i}{w_i + w_j}\right)$.*

Proof: The i th distance function is $f_i(\mathbf{q}) = w_i \|\mathbf{q} - \mathbf{p}_i\| + \alpha_i$. As such, for $y < \max(\alpha_i, \alpha_j)$, either $(f_i)_{\leq y} = \emptyset$, or $(f_j)_{\leq y} = \emptyset$, and $(f_i)_{\leq y} \cap (f_j)_{\leq y} = \emptyset$. For $y \geq \max(\alpha_i, \alpha_j)$, we have $f_i(\mathbf{q}) \leq y \iff w_i \|\mathbf{q} - \mathbf{p}_i\| + \alpha_i \leq y \iff \|\mathbf{q} - \mathbf{p}_i\| \leq \frac{y - \alpha_i}{w_i}$, which implies that $\mathbf{q} \in B\left(\mathbf{p}_i, \frac{y - \alpha_i}{w_i}\right)$; that is, we have $(f_i)_{\leq y} = B\left(\mathbf{p}_i, \frac{y - \alpha_i}{w_i}\right)$, and $(f_j)_{\leq y} = B\left(\mathbf{p}_j, \frac{y - \alpha_j}{w_j}\right)$.

Now, if $\mathbf{p}_i = \mathbf{p}_j$, then the distance between the two functions is the minimal value such that their sublevel sets are not empty, and this is $\max(\alpha_i, \alpha_j)$. In particular, $\frac{\alpha_i w_j + \alpha_j w_i}{w_i + w_j} \leq \max(\alpha_i, \alpha_j)$, and $\max\left(\alpha_i, \alpha_j, \|\mathbf{p}_i - \mathbf{p}_j\| \frac{w_i w_j}{w_i + w_j} + \frac{\alpha_i w_j + \alpha_j w_i}{w_i + w_j}\right) = \max(\alpha_i, \alpha_j)$, as desired.

If $\mathbf{p}_i \neq \mathbf{p}_j$ the sublevel sets intersect for the first time when the balls $\mathbf{B}\left(\mathbf{p}_i, \frac{y-\alpha_i}{w_i}\right)$ and $\mathbf{B}\left(\mathbf{p}_j, \frac{y-\alpha_j}{w_j}\right)$ touch at a point that belongs to the segment $\mathbf{p}_i\mathbf{p}_j$. Clearly then we have, $\|\mathbf{p}_i - \mathbf{p}_j\| = \frac{y-\alpha_i}{w_i} + \frac{y-\alpha_j}{w_j} \implies w_i w_j \|\mathbf{p}_i - \mathbf{p}_j\| = w_j(y - \alpha_i) + w_i(y - \alpha_j) \implies (w_i + w_j)y = w_i w_j \|\mathbf{p}_i - \mathbf{p}_j\| + w_j \alpha_i + w_i \alpha_j \implies y = \|\mathbf{p}_i - \mathbf{p}_j\| \frac{w_i w_j}{w_i + w_j} + \frac{\alpha_i w_j + \alpha_j w_i}{w_i + w_j}$. \blacksquare

Lemma 4.2. *Given $1 \leq i, j \leq n$, such that $w_i \leq w_j$. Suppose $y \geq \max(\alpha_i, \alpha_j)$. Then, for any $\delta \geq 0$, we have $(f_j)_{\leq y} \subseteq (f_i)_{\leq (1+\delta)y}$ if and only if $y \geq \frac{\|\mathbf{p}_i - \mathbf{p}_j\| + \alpha_i/w_i - \alpha_j/w_j}{(1+\delta)/w_i - 1/w_j}$.*

Proof: For $y \geq \max(\alpha_i, \alpha_j)$, we have that $(f_i)_{\leq y} = \mathbf{B}\left(\mathbf{p}_i, \frac{y-\alpha_i}{w_i}\right)$, and $(f_j)_{\leq y} = \mathbf{B}\left(\mathbf{p}_j, \frac{y-\alpha_j}{w_j}\right)$. If $\mathbf{p}_i = \mathbf{p}_j$, then for any y such that $\frac{(1+\delta)y-\alpha_i}{w_i} \geq \frac{y-\alpha_j}{w_j}$, we will have that $(f_j)_{\leq y} \subseteq (f_i)_{\leq (1+\delta)y}$. Clearly, this condition is also necessary. It is easy to verify that this is equivalent to the desired expression.

Consider the case $\mathbf{p}_i \neq \mathbf{p}_j$. For any $\mathbf{u} \in \mathbf{B}\left(\mathbf{p}_j, \frac{y-\alpha_j}{w_j}\right)$ we have $\|\mathbf{u} - \mathbf{p}_i\| \leq \|\mathbf{p}_i - \mathbf{p}_j\| + \|\mathbf{p}_j - \mathbf{u}\| \leq \|\mathbf{p}_i - \mathbf{p}_j\| + \frac{y-\alpha_j}{w_j}$, by the triangle inequality. Therefore, if $\frac{(1+\delta)y-\alpha_i}{w_i} \geq \|\mathbf{p}_i - \mathbf{p}_j\| + \frac{y-\alpha_j}{w_j}$, then $\mathbf{B}\left(\mathbf{p}_j, \frac{y-\alpha_j}{w_j}\right) \subseteq \mathbf{B}\left(\mathbf{p}_i, \frac{(1+\delta)y-\alpha_i}{w_i}\right)$. This is exactly the stated condition. Indeed, by rearrangement, $y((1+\delta)/w_i - 1/w_j) \geq \|\mathbf{p}_i - \mathbf{p}_j\| + \alpha_i/w_i - \alpha_j/w_j$.

As for the other direction, note that $\mathbf{B}\left(\mathbf{p}_j, \frac{y-\alpha_j}{w_j}\right)$ has a boundary point at distance $\frac{y-\alpha_j}{w_j}$ from \mathbf{p}_j on the directed line from \mathbf{p}_i to \mathbf{p}_j on the other side of \mathbf{p}_j as \mathbf{p}_i , while $\mathbf{B}\left(\mathbf{p}_i, \frac{(1+\delta)y-\alpha_i}{w_i}\right)$ has the intercept of $\frac{(1+\delta)y-\alpha_i}{w_i} - \|\mathbf{p}_i - \mathbf{p}_j\|$. For the condition to hold it must be true that $\frac{(1+\delta)y-\alpha_i}{w_i} - \|\mathbf{p}_i - \mathbf{p}_j\| \geq \frac{y-\alpha_j}{w_j}$, which is also the stated condition. \blacksquare

It is easy to see that compactness $(\text{P1})_{\text{p8}}$ and bounded growth $(\text{P2})_{\text{p8}}$ hold for the set of functions \mathcal{F} (for $(\text{P2})_{\text{p8}}$ we can take the growth function $\lambda_{(f_i)}(y) = (y - \alpha_i)/w_i$ for $y \geq \alpha_i$ and the growth constant ζ to be 2). The following lemma proves the sketch property $(\text{P3})_{\text{p8}}$.

Lemma 4.3. *For any $\mathcal{G} \subseteq \mathcal{F}$ and $\delta > 0$ there is a (δ, y_0) -sketch $\mathcal{H} \subseteq \mathcal{G}$, with $|\mathcal{H}| = 1$ and $y_0 = 3\ell_{\text{con}}(\mathcal{G})|\mathcal{G}|/\delta$.*

Proof: If $|\mathcal{G}| = 1$, set $\mathcal{H} = \mathcal{G}$. Otherwise, let $\ell = \ell_{\text{con}}(\mathcal{G})$ for brevity. Observe that $\ell \geq \max_{i: f_i \in \mathcal{G}} \alpha_i$, as otherwise some $(f_i)_{\leq \ell} = \emptyset$, and it cannot be part of a connected collection of sets. Let $|\mathcal{G}| = m \geq 2$, and let $\mathcal{G} = \{f_1, \dots, f_m\}$, and assume that $w_1 \leq w_i$, for $i = 1, \dots, m$. Set $\mathcal{H} = \{f_1\}$ – the function with the minimum associated weight. Since $\ell \geq \alpha_i$, we have $(f_i)_{\leq \ell}$ is the ball $\mathbf{B}\left(\mathbf{p}_i, \frac{\ell - \alpha_i}{w_i}\right)$, for $i = 1, \dots, m$. Since $\mathcal{G}_{\leq \ell}$ is connected, it must be true that, for a fixed $j, 2 \leq j \leq m$, there exist a sequence of distinct indices, $1 = i_1, i_2, \dots, i_{k-1}, i_k = j$, such that $\mathbf{B}\left(\mathbf{p}_{i_r}, \frac{\ell - \alpha_{i_r}}{w_{i_r}}\right) \cap \mathbf{B}\left(\mathbf{p}_{i_{r+1}}, \frac{\ell - \alpha_{i_{r+1}}}{w_{i_{r+1}}}\right) \neq \emptyset$, for $r = 1, \dots, k-1$. By Lemma 4.1, we have $\ell \geq \frac{\|\mathbf{p}_{i_r} - \mathbf{p}_{i_{r+1}}\| + \alpha_{i_r}/w_{i_r} + \alpha_{i_{r+1}}/w_{i_{r+1}}}{1/w_{i_r} + 1/w_{i_{r+1}}}$. Rearranging, $\|\mathbf{p}_{i_r} - \mathbf{p}_{i_{r+1}}\| \leq \ell \cdot \left(\frac{1}{w_{i_r}} + \frac{1}{w_{i_{r+1}}}\right) - \left(\frac{\alpha_{i_r}}{w_{i_r}} + \frac{\alpha_{i_{r+1}}}{w_{i_{r+1}}}\right) \leq \frac{2\ell}{w_1}$, as $w_1 \leq w_i$, for $i = 1, \dots, m$. It follows by the triangle inequality and the above, that $\|\mathbf{p}_{i_1} - \mathbf{p}_{i_k}\| \leq \sum_{r=1}^{k-1} \|\mathbf{p}_{i_r} - \mathbf{p}_{i_{r+1}}\| \leq 2(k-1)\ell/w_1 \leq 2m\ell/w_1$. Thus we have, $\|\mathbf{p}_1 - \mathbf{p}_j\| \leq 2m\ell/w_1$, for $j = 1, \dots, m$. Let $y_0 = 3\ell|\mathcal{G}|/\delta = 3m\ell/\delta$. Then, for $y \geq y_0$, we have that, $y \geq \frac{3m\ell}{\delta} = \frac{2m\ell/w_1 + m\ell/w_1}{\delta/w_1} \geq \frac{2m\ell/w_1 + \ell/w_1}{\delta/w_1}$, for $m \geq 2$. The above implies that, for $y \geq y_0$, we have $y \geq \frac{\|\mathbf{p}_1 - \mathbf{p}_j\| + \ell/w_1}{\delta/w_1} \geq \frac{\|\mathbf{p}_1 - \mathbf{p}_j\| + \alpha_1/w_1}{\delta/w_1}$, since $\ell \geq \alpha_1$. It follows that for $y \geq y_0$,

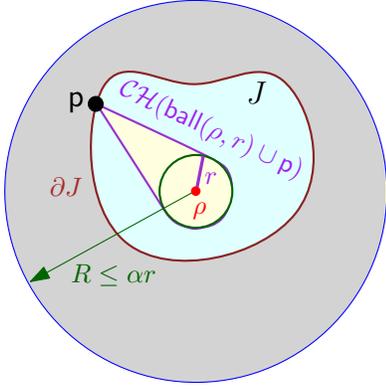


Figure 4.1: Being α -rounded fat.

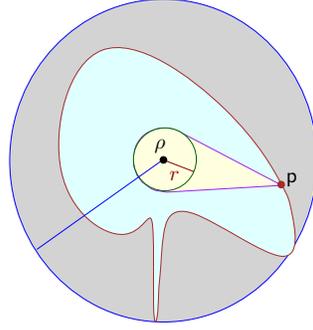


Figure 4.2: A α -fat star shaped region that is not α -rounded fat. Clearly, the gap between the rounded fatness parameter and the regular fatness parameter can be made to be arbitrarily large.

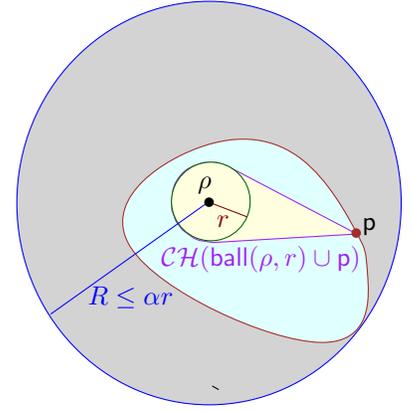


Figure 4.3: A α -fat convex body is α -rounded fat (for the same center point ρ).

$y \geq \frac{\|\mathbf{p}_1 - \mathbf{p}_j\| + \alpha_1/w_1}{\delta/w_1} \geq \frac{\|\mathbf{p}_1 - \mathbf{p}_j\| + \alpha_1/w_1 - \alpha_j/w_j}{\delta/w_1 + (1/w_1 - 1/w_j)} = \frac{\|\mathbf{p}_1 - \mathbf{p}_j\| + \alpha_1/w_1 - \alpha_j/w_j}{(1 + \delta)/w_1 - 1/w_j}$, as $w_1 \leq w_j$, for $j = 1, \dots, m$. Thus, by Lemma 4.2, $\mathbf{B}\left(\mathbf{p}_j, \frac{y - \alpha_j}{w_j}\right) \subseteq \mathbf{B}\left(\mathbf{p}_1, \frac{(1 + \delta)y - \alpha_1}{w_1}\right)$, for $y \geq y_0$, and therefore, by definition, \mathcal{H} is a (δ, y_0) -sketch for \mathcal{G} . \blacksquare

From the above, and the requirement that the sketch constant c_{sk} be at least $d + 1$ for our bounds to hold, it follows that we can choose the sketch constant $c_{\text{sk}} = \max(1, d + 1) = d + 1$. Moreover, as the graphs of the functions are cones, we can easily see that this function family satisfies the well-behavedness condition, see Definition 2.7_{p7}, and that we can choose the wellness constant $\gamma = 1$. Thus, we get Theorem 2.16_{p10}.

4.2. Scaling distance – generalized polytope distances

Let $J \subseteq \mathbb{R}^d$ be a compact set containing a “center” point ρ in its interior. Then J is *star shaped*, if for any point $\mathbf{v} \in J$, the entire segment $\rho\mathbf{v}$ is also in J . Naturally, any convex body J with any center ρ in the interior of J , is star shaped. The *t-scaling* of J with a center ρ , is the set, $tJ = \left\{ t(\mathbf{v} - \rho) + \rho \mid \mathbf{v} \in J \right\}$.

Given a star shaped object J with a center ρ , the *scaling distance* of a point \mathbf{q} from J is the minimum t , such that $\mathbf{p} \in tJ$, and let $f_J(\mathbf{q})$ denote this distance function. Note that, for any $y \geq 0$, the sublevel set $(f_J)_{\leq y}$, is the y -scaling of J , that is $(f_J)_{\leq y} = yJ$. Note that for a point $\mathbf{p} \in \mathbb{R}^d$, if we take $J = \mathbf{B}(\mathbf{p}, 1)$ with center \mathbf{p} , then $f_J(\mathbf{q}) = \|\mathbf{p} - \mathbf{q}\|$. That is, this distance notion is a strict extension of the Euclidean distance. Here, we assume that an object J contains the origin in its interior, and the origin is the designated center, unless otherwise stated.

Definition 4.4. A star shaped object $J \subseteq \mathbb{R}^d$ centered at ρ is *α -fat* if there is a number r such that, $\text{ball}(\rho, r) \subseteq J \subseteq \text{ball}(\rho, \alpha r)$.

Definition 4.5. Let J be a star shaped object centered at ρ . The object J is *α -rounded fat*, if there is a radius r such that, (i) $\text{ball}(\rho, r) \subseteq J \subseteq \text{ball}(\rho, \alpha r)$, and (ii) For every point \mathbf{p} in the boundary of J , the cone $\mathcal{CH}(\text{ball}(\rho, r) \cup \mathbf{p})$, lies within J , see Figure 4.1.

By definition, any α -rounded fat object is also α -fat. However, it is not true that a α -fat object is necessarily rounded fat, see Figure 4.2. The following useful result is easy to see, also see Figure 4.3 for an illustration.

Lemma 4.6. *An object J that is a α -fat and convex is also α -rounded fat.*

Given a set $\mathcal{J} = \{J_1, J_2, \dots, J_n\}$ 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 $f_i = f_{J_i}$. We assume that the boundary of each object J_i , has constant complexity.

We next argue that \mathcal{F} complies with the framework of Section 2.4. Using standard techniques, we can compute the quantities required in conditions (C1)–(C3)_{p8}, as well as the diameter of the sublevel sets, as $\text{diam}(yJ_i) = y\text{diam}(J_i)$. Also, trivially we have that condition (P1)_{p8} is satisfied, as the sublevel sets are dilations of the J_i , and are thus compact by definition. The next few lemmas establish that both bounded growth (P2)_{p8}, and the sketch property (P3)_{p8}, are also true, if the objects are also α -rounded fat, for some constant α .

Lemma 4.7. *Given $\alpha > 0$, and an object J that is α -rounded fat. Then, for any $c \geq 2\alpha$, $y \geq 0$, and $\varepsilon > 0$, we have that $yJ \oplus \mathbf{B}(0, (\varepsilon/c)\text{diam}(yJ)) \subseteq (1 + \varepsilon)yJ$; that is, $(f_J)_{\leq y} \oplus \mathbf{B}\left(0, (\varepsilon/c)\text{diam}\left((f_J)_{\leq y}\right)\right) \subseteq (f_J)_{\leq (1+\varepsilon)y}$.*

Proof: Since $(f_J)_{\leq y} = yJ$ it is enough to show that $yJ \oplus \mathbf{B}(0, (\varepsilon/c)\text{diam}(yJ)) \subseteq (1 + \varepsilon)yJ$. Let r be the radius guaranteed by Definition 4.5 for J . Clearly $\text{diam}(yJ) = y\text{diam}(J) \leq 2y\alpha r$. Now, for every $\mathbf{p} \in \partial yJ$, we have that $\mathbf{p} + \mathbf{B}(0, (\varepsilon/c)\text{diam}(yJ)) \subseteq (1 + \varepsilon)yJ$. It is sufficient to show that $\mathbf{B}(\mathbf{p}, (2\varepsilon y\alpha r/c)) \subseteq (1 + \varepsilon)yJ$, for $\mathbf{p} \in \partial yJ$. Clearly $\mathbf{p}' = (1 + \varepsilon)\mathbf{p} \in \partial(1 + \varepsilon)yJ$. Since the cone, $\mathcal{CH}(\mathbf{B}(\rho, (1 + \varepsilon)yr) \cup \mathbf{p}')$ is in $(1 + \varepsilon)yJ$, it is clear that the ball of radius, $x = \|\mathbf{p}' - \mathbf{p}\| \frac{(1 + \varepsilon)yr}{\|\mathbf{p}'\|} = \|\mathbf{p}' - \mathbf{p}\| \frac{yr}{\|\mathbf{p}\|}$ is completely within $(1 + \varepsilon)yJ$, see Figure 4.4. Now, $\|\mathbf{p} - \mathbf{p}'\| = \varepsilon \|\mathbf{p}\|$, so $x = \varepsilon yr$. For $c \geq 2\alpha$, the result follows. ■

By the above lemma, we can take the growth function $\lambda_{f_{J_i}}(y) = \text{diam}\left((f_{J_i})_{\leq y}\right) = y\text{diam}(J_i)$, and the growth constant, see (P2)_{p8}, for the set of functions f_{J_i} , to be $\zeta = c = 2\alpha$. If the object J is α -fat, but not α' -rounded fat for any constant $\alpha' > 0$, then it may be that its scaling distance function grows arbitrarily quickly, and thus is not a valid distance function for our framework, see Figure 4.5. It is not hard to see that Lemma 4.7 implies that bounded growth (P2)_{p8} is satisfied for all the functions f_1, \dots, f_n , when the objects under consideration J_1, \dots, J_n , are α -rounded fat.

4.2.1. There is a small sketch and it can be computed quickly

To show that a small sketch exists (i.e., (P3)_{p8} holds) is slightly harder and is tackled next.

Lemma 4.8. *Let $\mathcal{J} = \{J_1, J_2, \dots, J_n\}$ be a set of n α -rounded fat objects centered at the origin, where $\alpha \geq 1$ is a fixed constant. Then, for any $\delta > 0$, there is a subset $\mathcal{I} \subseteq \{1, 2, \dots, n\}$ with $|\mathcal{I}| = O(\delta^{-d})$, such that for all $y \geq 0$, we have $\bigcup_{i \in \{1, 2, \dots, n\}} yJ_i \subseteq \bigcup_{i \in \mathcal{I}} (1 + \delta)yJ_i$. Moreover, for every $i \in \mathcal{I}$, we have that $\text{diam}(J_i) = \Omega(\max_i \text{diam}(J_i))$.*

Proof: Clearly it is sufficient to show this for $y = 1$. For $i = 1, \dots, n$, let r_i be the radius of the “fatness” ball contained in J_i , see Definition 4.5_{p22}. We assume without loss of generality that r_1 is the maximum

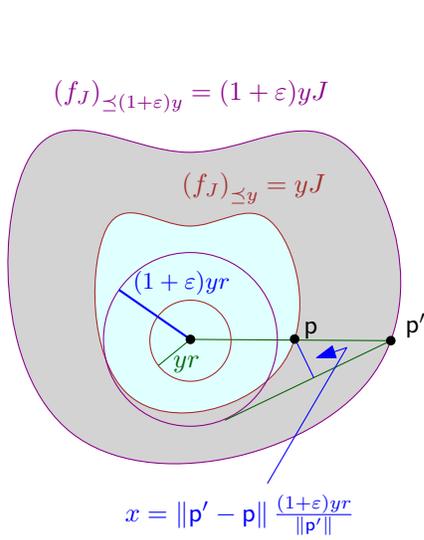


Figure 4.4: The $(1 + \varepsilon)$ expansion of yJ contains $\mathbf{B}(\mathbf{p}, x)$.

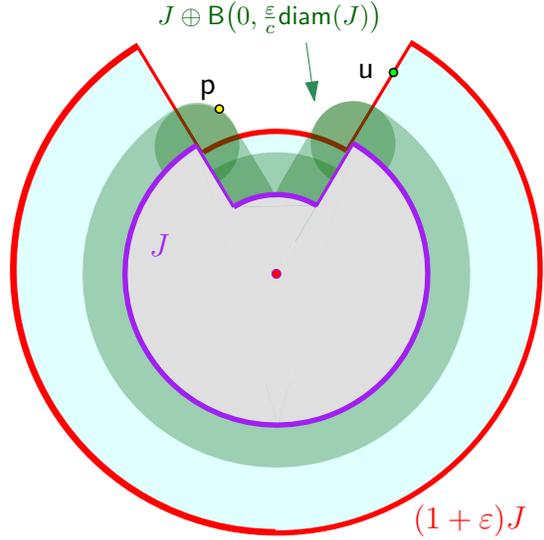


Figure 4.5: The object J is α -fat but not α' -rounded fat for any $\alpha' > 0$. In particular, the point \mathbf{p} is in $J \oplus \mathbf{B}(0, (\varepsilon/c)\text{diam}(J))$ but not in $(1 + \varepsilon)J$, and the scaling distance function is not well defined at \mathbf{u} .

of all the r_i . Let $\mathcal{K} = \left\{ i \mid \alpha r_i \geq r_1 \right\}$ be the set of indices of relatively large objects (the other objects are sufficiently small and are contained in J_1). Observe that \mathcal{K} is not empty as $1 \in \mathcal{K}$. We have that

$$\text{diam}(J_i) \geq 2r_i \geq \frac{2}{\alpha} r_1 \geq \frac{1}{\alpha^2} \cdot 2\alpha r_1 \geq \frac{1}{\alpha^2} \max_{1 \leq i \leq n} \text{diam}(J_i),$$

for any $i \in \mathcal{K}$, since all the objects are centered at the origin.

Clearly, $\bigcup_{i=1}^n J_i \subseteq \bigcup_{i=1}^n \mathbf{B}(0, \alpha r_i) \subseteq \mathbf{B}(0, \alpha r_1)$. We tile the ball $\mathbf{B}(0, \alpha r_1)$ with cubes of diameter $\delta \alpha r_1 / c'$ where $c' = c\alpha^2/2 = \alpha^3$. Let \mathcal{C} denote the set of these cubes, and observe that $|\mathcal{C}| = O(\delta^{-d})$. For every cube $\mathbf{c} \in \mathcal{C}$, if it intersects $\mathbf{X} = \bigcup_{i \in \mathcal{K}} J_i$, then we add the index of one of the objects of \mathcal{K} intersecting \mathbf{c} to \mathcal{I} , and add \mathbf{c} to \mathcal{A} .

Now, $\bigcup_{i=1}^n J_i \subseteq \bigcup_{\mathbf{c} \in \mathcal{A}} \mathbf{c}$, as clearly we have that,

$$\bigcup_{i=1}^n J_i = \mathbf{B}(0, \alpha r_1) \cap \left(\bigcup_{i=1}^n J_i \right) \subseteq \left(\bigcup_{\mathbf{c} \in \mathcal{C}} \mathbf{c} \right) \cap \left(\bigcup_{i=1}^n J_i \right) \subseteq \bigcup_{\mathbf{c} \in \mathcal{A}} \mathbf{c}.$$

Observe that $|\mathcal{I}| = |\mathcal{A}| \leq |\mathcal{C}| = O(\delta^{-d})$. Next, we show that, $\bigcup_{\mathbf{c} \in \mathcal{A}} \mathbf{c} \subseteq \bigcup_{i \in \mathcal{I}} (1 + \delta) J_i$. Suppose $\mathbf{c} \in \mathcal{A}$ and let $i \in \mathcal{I}$ be an index such that $\mathbf{c} \cap J_i \neq \emptyset$. Since $\text{diam}(\mathbf{c}) \leq \delta \alpha r_1 / c'$, $\mathbf{c} \subseteq J_i \oplus \mathbf{B}(0, \delta \alpha r_1 / c')$. Now, $\frac{\delta \alpha r_1}{c'} \leq \frac{\delta \alpha^2 r_i}{c'} \leq \frac{\delta \alpha^2 \text{diam}(J_i)}{2c'} \leq \frac{\delta \text{diam}(J_i)}{c}$, where $c = 2\alpha$ is the constant from Lemma 4.7. Then, by Lemma 4.7, we have

$$\mathbf{c} \subseteq J_i \oplus \mathbf{B}(0, \delta \alpha r_1 / c') \subseteq J_i \oplus \mathbf{B}(0, \delta \text{diam}(J_i) / c) \subseteq (1 + \delta) J_i,$$

which implies the claim. ■

Lemma 4.9. *Let $\alpha \geq 1$ be a constant, J an α -rounded fat object, $\delta > 0$, and let \mathbf{u} be a point in \mathbb{R}^d , with $\|\mathbf{u}\| \leq \delta \text{diam}(J) / c$, where $c = 2\alpha$. Then $J + \mathbf{u} \subseteq (1 + \delta) J$.*

Proof: We have, $J + \mathbf{u} \subseteq J \oplus \mathbf{B}(0, \|\mathbf{u}\|) \subseteq J \oplus \mathbf{B}(0, \delta \text{diam}(J) / c)$ as $\|\mathbf{u}\| \leq \delta \text{diam}(J) / c$. Now, the claim follows by Lemma 4.7. ■

Lemma 4.10. *Let $\mathcal{J} = \{J_1, \dots, J_n\}$ be a set of n α -rounded fat object in \mathbb{R}^d , let \mathcal{F} be the set of scaling distance functions of these objects, and let \mathbf{P} be the set of offset points (i.e., centers) of the objects of \mathcal{J} . Then $\ell_{\text{con}}(\mathcal{F}) \geq \text{diam}(\mathbf{P})/2n\alpha r$, where $r = \max_i r_i$, and r_i is the inner radius of J_i , see Definition 4.5_{p22}.*

Proof: Let $\mathcal{F} = \{f_i = f_{J_i} \mid i = 1, \dots, n\}$ and $\mathbf{P} = \{\mathbf{p}_1, \dots, \mathbf{p}_n\}$. Here, for $i = 1, \dots, n$, \mathbf{p}_i is the center of J_i . The claim is trivially true if $\text{diam}(\mathbf{P}) = 0$, i.e., all the points \mathbf{p}_i are the same. Let $\ell = \ell_{\text{con}}(\mathcal{F})$. As $(f_i)_{\leq \ell} = \ell J_i$, where the scaling for object J_i is done around its center \mathbf{p}_i , it follows that $\ell \mathcal{J} = \{\ell J_i \mid J_i \in \mathcal{F}\}$ is connected, see Definition 2.6. Since $\ell J_i \subseteq \mathbf{B}(\mathbf{p}_i, \ell \alpha r_i) \subseteq \mathbf{B}(\mathbf{p}_i, \ell \alpha r)$, it follows that $\mathcal{B} = \{\mathbf{B}(\mathbf{p}_i, \ell \alpha r) \mid i = 1, \dots, n\}$ is also connected. Let $\mathbf{u}, \mathbf{v} \in \mathbf{P}$ be such that $\|\mathbf{u} - \mathbf{v}\| = \text{diam}(\mathbf{P})$. There is a sequence of distinct $i_1, \dots, i_k \in \{1, \dots, n\}$, such that $\mathbf{u} = \mathbf{p}_{i_1}, \mathbf{v} = \mathbf{p}_{i_k}$, and $\mathbf{B}(\mathbf{p}_{i_t}, \ell \alpha r) \cap \mathbf{B}(\mathbf{p}_{i_{t+1}}, \ell \alpha r) \neq \emptyset$, and thus $\|\mathbf{p}_{i_t} - \mathbf{p}_{i_{t+1}}\| \leq 2\ell \alpha r$, for $t = 1, \dots, k-1$. By the triangle inequality,

$$\text{diam}(\mathbf{P}) = \|\mathbf{u} - \mathbf{v}\| = \|\mathbf{p}_{i_1} - \mathbf{p}_{i_k}\| \leq \sum_{t=1}^{k-1} \|\mathbf{p}_{i_t} - \mathbf{p}_{i_{t+1}}\| \leq \sum_{t=1}^{k-1} 2\ell \alpha r = 2(k-1)\ell \alpha r \leq 2n\ell \alpha r, \quad \blacksquare$$

We can now show that condition (P3)_{p8} holds for the f_{J_i} .

Lemma 4.11. *Consider the setting of Lemma 4.10. Given $\delta > 0$, there is a index set $\mathcal{I} \subseteq \{1, \dots, n\}$, with $|\mathcal{I}| = O(\delta^{-d})$, and $y_0 = O(\ell \cdot n/\delta)$, such that the functions $\{f_j \mid j \in \mathcal{I}\}$, form a (δ, y_0) -sketch, where $\ell = \ell_{\text{con}}(\mathcal{F})$.*

Proof: We provide a sketch of the proof, as the details are easy but tedious. Consider the set of objects $J_{ij} = J_i + \mathbf{p}_j - \mathbf{p}_i$, for each pair (i, j) with $1 \leq i, j \leq n$. Clearly, J_{ij} is J_i translated, so that it is centered at \mathbf{p}_j . By Lemma 4.8 there is an index set $\mathcal{I} \subseteq \{1, \dots, n\}$, with $|\mathcal{I}| = O(\delta^{-d})$, such that for all y and any fixed j with $1 \leq j \leq n$, we have that $\bigcup_{i=1}^n y J_{ij} \subseteq \bigcup_{i \in \mathcal{I}} (1 + \delta/4)y J_{ij}$. Let r_i denote the radius of the ball for J_i from Definition 4.5, and let $r = \max_i r_i$. By Lemma 4.10, we have that, $\ell \geq \text{diam}(\mathbf{P})/(2n\alpha r)$. Lemma 4.8 finds a \mathcal{I} such that for all $i \in \mathcal{I}$, $r_i \geq \Omega(r)$. A translated copy $J_{ij} = J_i + \mathbf{p}_j - \mathbf{p}_i$, is a translation of J_i by a vector $\mathbf{u} = \mathbf{p}_j - \mathbf{p}_i$. Now for any $y \geq 0$,

$$\bigcup_{i=1}^n y J_i \subseteq \bigcup_{j=1}^n \bigcup_{i=1}^n y J_{ij} \subseteq \bigcup_{j=1}^n \bigcup_{i \in \mathcal{I}} (1 + \delta/4)y J_{ij} = \bigcup_{j=1}^n \bigcup_{i \in \mathcal{I}} ((1 + \delta/4)y J_i + (\mathbf{p}_j - \mathbf{p}_i)).$$

As $\ell \geq \text{diam}(\mathbf{P})/(2n\alpha r)$, there is a $y_0 = O(\ell n/\delta)$ such that,

$$\|\mathbf{p}_j - \mathbf{p}_i\| \leq \delta \text{diam}(y_0 J_i)/4c \leq \delta \text{diam}((1 + \delta/4)y_0 J_i)/4c,$$

for all $1 \leq i, j \leq n$, where $c = 2\alpha$. Thus using Lemma 4.9 for $y = y_0$ on a $(1 + \delta/4)$ expansion of the J_i we have that,

$$\bigcup_{=1}^n \bigcup_{i \in \mathcal{I}} (1 + \delta/4)y J_i + (\mathbf{p}_j - \mathbf{p}_i) \subseteq \bigcup_{j=1}^n \bigcup_{i \in \mathcal{I}} (1 + \delta/4)^2 y J_i \subseteq \bigcup_{j=1}^n \bigcup_{i \in \mathcal{I}} (1 + \delta)y J_i = \bigcup_{i \in \mathcal{I}} (1 + \delta)y J_i.$$

It is easy to verify that this also holds for any $y \geq y_0$. As such by definition, $\{f_i \mid i \in \mathcal{I}\}$, is a (δ, y_0) -sketch. \blacksquare

4.2.2. The result

We conclude that for α -rounded fat objects, the scaling distance functions they define, falls under our framework. From the above, and the requirement that the sketch constant c_{sk} be at least $d + 1$ for our bounds to hold, it follows that we can choose the sketch constant $c_{\text{sk}} = \max(d, d + 1) = d + 1$, and this concludes the proof of Theorem 2.17_{p10}.

Note that the result in Theorem 2.17 covers any symmetric convex metric. Indeed, given a convex symmetric shape C centered at the origin, and with constant boundary complexity, the distance it induces for any pair of points $\mathbf{p}, \mathbf{u} \in \mathbb{R}^d$, is the scaling distance of C centered at \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 2.17 applies. Of course, Theorem 2.17 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.

4.3. 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 $F_S(\mathbf{q}) = \max_{s \in S} \|\mathbf{q} - s\|$; 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 $F_S(\cdot)$ 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 (for constant dimensions), see [Har99] – one can compute, in linear time, a subset $U \subseteq S$, of size $O(\varepsilon^{-d} \log \varepsilon^{-1})$, such that $(1 - \varepsilon)F_S(\mathbf{q}) \leq F_U(\mathbf{q}) \leq F_S(\mathbf{q})$.

Given n sets of points P_1, \dots, P_n in \mathbb{R}^d , we are interested in the distance function $F(\mathbf{q}) = \min_i F_i(\mathbf{q})$, where $F_i(\mathbf{q}) = F_{P_i}(\mathbf{q})$. This quantity arises naturally when one tries to model uncertainty; indeed, let P_i be the set of possible locations of the i th *uncertain* point (i.e., the location of the i th point is chosen randomly, somehow, from the set P_i). Thus, $F_i(\mathbf{q})$ is the worst case distance to the i th point, and $F(\mathbf{q})$ 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 $F(\cdot)$ as the *nearest furthest-neighbor* Voronoi diagram, and we are interested in its approximation.

4.4. Satisfaction of conditions

Observation 4.12. *We have that $(F_i)_{\leq y} = \bigcap_{u \in P_i} B(u, y)$, and $\text{diam}\left((F_i)_{\leq y}\right) \leq 2y$.*

Given the above observation, it is easy to see that Condition (P1)_{p8} is true, as $(F_i)_{\leq y}$ is a finite intersection of compact sets. The following Lemma shows that Condition (P2)_{p8} is also true, by letting the growth function $\lambda_{(F_i)}(y) = y$. Since $y \geq \text{diam}\left((F_i)_{\leq y}\right)/2$ by Observation 4.12, it follows that we can choose the growth constant ζ to be 2.

Lemma 4.13. *For any i with $1 \leq i \leq n$, if $(F_i)_{\leq y} \neq \emptyset$, we have $(F_i)_{\leq y} \oplus B(0, \varepsilon y) \subseteq (F_i)_{\leq (1+\varepsilon)y}$.*

Proof: Consider any point \mathbf{q} in $(F_i)_{\leq y} \oplus B(0, \varepsilon y)$. Now, one can easily verify that $B(\mathbf{q}, (1 + \varepsilon)y) \supseteq P_i$ by the triangle inequality, and so $\mathbf{q} \in (F_i)_{\leq (1+\varepsilon)y}$. ■

Lemma 4.14 (Existence of a sketch). *Let $\mathcal{G} \subseteq \{F_1, \dots, F_n\}$, denote a set of functions as above. Then, given any $\delta > 0$, there is a subset $\mathcal{H} \subseteq \mathcal{G}$ with $|\mathcal{H}| = 1$, and a y_0 with $y_0 = O(\ell_{\text{con}}(\mathcal{G}) |\mathcal{G}| / \delta)$, such that \mathcal{H} is a (δ, y_0) -sketch for \mathcal{G} .*

Proof: Let $\mathcal{G} = \{F_1, \dots, F_m\}$, where $m = |\mathcal{G}|$, and $z = \ell_{\text{con}}(\mathcal{G})$. Now, $(F_i)_{\leq z}$ for $i = 1, \dots, m$, are all connected, and by Observation 4.12 $\text{diam}((F_i)_{\leq z}) \leq 2z$. Thus, $\forall \mathbf{u} \in P_j, \forall \mathbf{v} \in P_k$, there are points $\mathbf{u}' \in (F_j)_{\leq z}$, and $\mathbf{v}' \in (F_k)_{\leq z}$, such that $\|\mathbf{u} - \mathbf{u}'\| \leq z$, $\|\mathbf{v} - \mathbf{v}'\| \leq z$ (by definition of the function F_j and F_k respectively), and $\|\mathbf{u}' - \mathbf{v}'\| \leq 2mz$, by the bound on the diameter of the sublevel sets and the condition of being connected, which is the same as the intersection graph of the sets being connected. It follows, by the triangle inequality, that $\|\mathbf{u} - \mathbf{v}\| \leq 2(m+1)z \leq 4mz$; namely, $\text{diam}(P_1 \cup \dots \cup P_m) \leq 4mz$. Let \mathcal{H} be a set containing an arbitrary function from \mathcal{G} say, $\mathcal{H} = \{F_1\}$.

It is easy to see, that for $y \geq y_0 = 4mz/\delta$, and for all $\mathbf{u}, \mathbf{v} \in P_1 \cup \dots \cup P_m$, we have that $B(\mathbf{v}, y) \subseteq B(\mathbf{u}, (1+\delta)y)$. Thus, for any i , $1 \leq i \leq m$, we have that $(F_i)_{\leq y} = \bigcap_{\mathbf{v} \in P_i} B(\mathbf{v}, y) \subseteq B(\mathbf{u}, (1+\delta)y)$, for all $\mathbf{u} \in P_1$ and $y \geq y_0$. As such, $(F_i)_{\leq y} \subseteq \bigcap_{\mathbf{u} \in P_1} B(\mathbf{u}, (1+\delta)y) = (F_1)_{\leq (1+\delta)y}$, and the result follows. ■

Remark 4.15. It is not hard, if somewhat tedious, to verify that conditions (C1)–(C3)_{p8} hold. We need to use a data-structure for approximate furthest-neighbor queries, as well as coresets for approximately computing the MEB of the point sets. Significantly, the center of the MEB is not perturbed too much when using coresets – see Lemma A.1 in Appendix A.

4.4.1. The result

From the above, and the requirement that the sketch constant c_{sk} be at least $d+1$ for our bounds to hold, it follows that for the functions under consideration, we can choose the sketch constant $c_{\text{sk}} = \max(d, d+1) = d+1$.

Proof of Theorem 2.18_{p11}: We only need to show how to get the improved space and query time – i.e., how do we get rid of the total number of points m in the 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 $F_{S_i}(\mathbf{q}) \leq F_{P_i}(\mathbf{q}) \leq (1+\varepsilon/4)F_{S_i}(\mathbf{q})$. Such a subset can be computed in $O(|P_i|)$ time, see [Har99]. We thus perform this transformation for each one of the uncertain point sets P_1, \dots, P_n , which reduces the input size to $O(n/\varepsilon^d \log(1/\varepsilon))$. It is easy to verify that all the other operations required to construct the data-structure can be done efficiently – for example, computing the sketch requires approximating the diameter up to a constant factor, which can easily be done in linear time. We now apply our main result to the distance functions induced by the reduced sets S_1, \dots, S_n . ■

5. 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 such functions, we presented a new data-structure, with near linear space and preprocessing time, so that it 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, multiplicative weighted distance functions with additive offsets. 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 new AVD has near linear complexity). The framework also works for Minkowski metrics of fat convex bodies, and nearest furthest-neighbor. However, it seems that 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 [AM02], should work in the most general settings, so reducing the logarithmic dependency seems quite interesting. Specifically, can the Arya and Malamatos construction [AM02] 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 dependable family amenable to our framework? Specifically, consider a polynomial $\tau(x)$ that is positive for all $x \geq 0$. Given a point \mathbf{u} , we associate the distance function $f(\mathbf{q}) = \tau(\|\mathbf{q} - \mathbf{u}\|)$ with \mathbf{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 under our framework.)

Acknowledgments. The authors thank the anonymous referees for their insightful and useful comments on the paper.

References

- [AAH⁺13] 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.* (PODS), pages 115–126, 2013.
- [AESZ12] P. K. Agarwal, A. Efrat, S. Sankararaman, and W. Zhang. Nearest-neighbor searching under uncertainty. In *Proc. 31st ACM Sympos. Principles Database Syst.* (PODS), pages 225–236, 2012.
- [Agg09] C. Aggarwal. *Managing and Mining Uncertain Data*. Springer, 2009.
- [AI08] A. Andoni and P. Indyk. Near-optimal hashing algorithms for approximate nearest neighbor in high dimensions. *Commun. ACM*, 51(1):117–122, 2008.
- [AM93] P. K. Agarwal and J. Matoušek. Ray shooting and parametric search. *SIAM J. Comput.*, 22:540–570, 1993.
- [AM02] S. Arya and T. Malamatos. Linear-size approximate Voronoi diagrams. In *Proc. 13th ACM-SIAM Sympos. Discrete Algs.* (SODA), pages 147–155, 2002.
- [AMM09] S. Arya, T. Malamatos, and D. M. Mount. Space-time tradeoffs for approximate nearest neighbor searching. *J. Assoc. Comput. Mach.*, 57(1):1–54, 2009.
- [AMN⁺98] 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.*, 45(6):891–923, 1998.

- [BCKO08] M. de Berg, O. Cheong, M. van Kreveld, and M. H. Overmars. *Computational Geometry: Algorithms and Applications*. Springer-Verlag, 3rd edition, 2008.
- [BPR06] S. Basu, R. Pollack, and M. F. Roy. *Algorithms in Real Algebraic Geometry*. Algorithms and Computation in Mathematics. Springer, 2006.
- [Cha10] T. M. Chan. Optimal partition trees. In *Proc. 26th Annu. Sympos. Comput. Geom. (SoCG)*, pages 1–10, 2010.
- [CK95] 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.*, 42:67–90, 1995.
- [Cla88] K. L. Clarkson. A randomized algorithm for closest-point queries. *SIAM J. Comput.*, 17:830–847, 1988.
- [Cla06] K. L. Clarkson. Nearest-neighbor searching and metric space dimensions. In G. Shakhnarovich, T. Darrell, and P. Indyk, editors, *Nearest-Neighbor Methods for Learning and Vision: Theory and Practice*, pages 15–59. MIT Press, 2006.
- [DRS09] N. N. Dalvi, C. Ré, and D. Suciu. Probabilistic databases: Diamonds in the dirt. *Commun. ACM*, 52(7):86–94, 2009.
- [Eri96] J. Erickson. New lower bounds for Hopcroft’s problem. *Discrete Comput. Geom.*, 16:389–418, 1996.
- [Har99] S. Har-Peled. Constructing approximate shortest path maps in three dimensions. *SIAM J. Comput.*, 28(4):1182–1197, 1999.
- [Har01] S. Har-Peled. A replacement for Voronoi diagrams of near linear size. In *Proc. 42nd Annu. IEEE Sympos. Found. Comput. Sci. (FOCS)*, pages 94–103, 2001.
- [Har11] S. Har-Peled. *Geometric Approximation Algorithms*, volume 173 of *Mathematical Surveys and Monographs*. Amer. Math. Soc., 2011.
- [HIM12] S. Har-Peled, P. Indyk, and R. Motwani. Approximate nearest neighbors: Towards removing the curse of dimensionality. *Theory Comput.*, 8:321–350, 2012. Special issue in honor of Rajeev Motwani.
- [HK13] S. Har-Peled and N. Kumar. Approximating minimization diagrams and generalized proximity search. In *Proc. 54th Annu. IEEE Sympos. Found. Comput. Sci. (FOCS)*, pages 717–726, 2013.
- [IM98] P. Indyk and R. Motwani. Approximate nearest neighbors: Towards removing the curse of dimensionality. In *Proc. 30th Annu. ACM Sympos. Theory Comput. (STOC)*, pages 604–613, 1998.
- [Mat92] J. Matoušek. Efficient partition trees. *Discrete Comput. Geom.*, 8:315–334, 1992.
- [Mei93] S. Meiser. Point location in arrangements of hyperplanes. *Inform. Comput.*, 106:286–303, 1993.

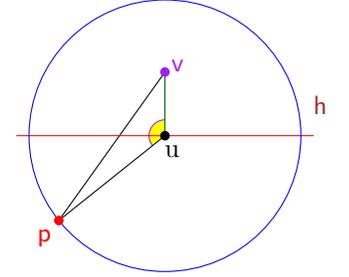
- [MNP06] R. Motwani, A. Naor, and R. Panigrahi. Lower bounds on locality sensitive hashing. In *Proc. 22nd Annu. Sympos. Comput. Geom. (SoCG)*, pages 154–157, 2006.
- [SA95] M. Sharir and P. K. Agarwal. *Davenport-Schinzel Sequences and Their Geometric Applications*. Cambridge University Press, New York, 1995.

A. Bounding the size of intersection of balls of the same radius

Lemma A.1. *Let $\text{ball}(u, z)$ be the MEB (minimum enclosing ball) of a set of points $P \subseteq \mathbb{R}^d$. Then, for any $\delta \in (0, 1)$, we have $\text{ball}(u, \delta z) \subseteq \bigcap_{p \in P} B(p, (1 + \delta)z) \subseteq \text{ball}\left(u, \sqrt{4\delta + 2\delta^2}z\right)$.*

Proof: There are affinely independent points from P , on the surface of the MEB, such that u lies in their convex hull (otherwise the MEB is not minimal). Thus, there is a set $S = \{p_1, \dots, p_k\} \subseteq P$, such that (i) $\forall p \in S \quad \|p - u\| = z$, and (ii) $u \in \mathcal{CH}(S)$. Thus, any closed halfspace h^+ such that its boundary passes through u must contain at least one of the points of S .

Consider any point $v \in \bigcap_{p \in S} B(p, (1 + \delta)z)$. Let h be the hyperplane passing through u , and orthogonal to $u - v$, and let h^+ be the closed halfspace having h on its boundary that does not contain v . Then, by the above observation, there must be a point $p \in S$, such that $p \in h^+$. Now, by the law of cosines, we have that



$$\begin{aligned} (1 + \delta)^2 z^2 &\geq \|v - p\|^2 = \|v - u\|^2 + \|u - p\|^2 - 2 \|v - u\| \|u - p\| \cos \angle vup \\ &\geq \|v - u\|^2 + \|u - p\|^2 = \|v - u\|^2 + z^2, \end{aligned}$$

since $\angle vup \geq \pi/2$. Rearranging, we have $(2\delta + \delta^2)z^2 \geq \|v - u\|^2$, thus implying $\sqrt{2\delta + \delta^2}z \geq \|v - u\|$. We conclude that $\bigcap_{p \in P} B(p, (1 + \delta)z) \subseteq \bigcap_{p \in S} B(p, (1 + \delta)z) \subseteq \text{ball}\left(u, \sqrt{2\delta + \delta^2}z\right)$. ■

B. Basic properties of the functions

Lemma B.1. *Let \mathcal{F} be a set of functions that satisfy the compactness $(P1)_{p8}$ and bounded growth $(P2)_{p8}$ conditions. Then, for any $f \in \mathcal{F}$, either $f_{\leq 0} = \emptyset$ or $f_{\leq 0}$ consists of a single point.*

Proof: If $f_{\leq 0}$ contains at least two points, then by compactness it contains two distinct points x and y , such that $\|x - y\| = \text{diam}(f_{\leq 0}) > 0$. By bounded growth $(P2)_{p8}$, with $\varepsilon = 1$, it follows

$$f_{\leq 0} \subseteq f_{\leq 0} \oplus B\left(0, \frac{\|x - y\|}{\zeta}\right) \subseteq f_{\leq 0} \oplus B(0, \lambda_f(0)) \subseteq f_{\leq 0},$$

as $\lambda_f(0) \geq \text{diam}(f_{\leq 0})/\zeta = \|x - y\|/\zeta$. Thus, $f_{\leq 0} \oplus B\left(0, \frac{\|x - y\|}{\zeta}\right) = f_{\leq 0}$. Clearly $y \oplus B\left(0, \frac{\|x - y\|}{\zeta}\right)$ contains a point y' such that $\|x - y'\| > \|x - y\|$. Namely, $\text{diam}(f_{\leq 0}) > \text{diam}(f_{\leq 0})$, which is a contradiction. ■

We assume that that $\text{d}(f, g) > 0$ for $f \neq g$ (this can be enforced using symbolic perturbations). Similarly, we also assume that the distances $\text{d}(f, g)$ are distinct for all pairs of functions.

Observation B.2. *If $\ell_{\text{con}}(\mathcal{F}) = 0$ for any non-empty set \mathcal{F} , then $|\mathcal{F}| = 1$.*

Lemma B.3. Let $f \in \mathcal{F}$ and $y \geq 0$. Suppose $\mathbf{u}, \mathbf{v} \in f_{\leq y}$. Then, $\mathbf{uv} \subseteq \mathcal{F}_{\leq (1+\zeta/2)y}$, where \mathbf{uv} denotes the segment joining \mathbf{u} to \mathbf{v} .

Proof: If $\mathbf{u} = \mathbf{v}$, the claim is obvious. Using bounded growth $(\text{P2})_{\text{p8}}$ with $\varepsilon = \zeta/2$, and the inequality $\lambda_f(y) \geq \text{diam}(f_{\leq y})/\zeta$, it follows that $f_{\leq y} \oplus \mathbf{B}(0, \text{diam}(f_{\leq y})/2) \subseteq f_{\leq (1+\zeta/2)y}$. Thus, $\mathbf{u} \oplus \mathbf{B}(0, \text{diam}(f_{\leq y})/2) \subseteq f_{\leq (1+\zeta/2)y}$ and $\mathbf{B}(\mathbf{v}, \text{diam}(f_{\leq y})/2) \subseteq f_{\leq (1+\zeta/2)y}$. Since $\|\mathbf{u} - \mathbf{v}\| \leq \text{diam}(f_{\leq y})$, it follows that the entire segment \mathbf{uv} is in $f_{\leq (1+\zeta/2)y}$. ■

Observation B.4. Let $A_1, \dots, A_m \subseteq \mathbb{R}^d$, be compact connected sets. Let \mathbf{uv} be any segment. Suppose that $\mathbf{uv} \cap A_i \neq \emptyset$, for $i = 1, \dots, k$, and $\mathbf{uv} \subseteq \bigcup_{j=1}^k A_j$. Then, the set $\{A_1, \dots, A_k\}$ is connected, see Definition 2.6_{p7}.

Lemma B.5. Given $\mathcal{G} \subseteq \mathcal{F}$, $\delta \geq 0$ and $y \geq 0$, such that \mathcal{G} is a (δ, y) -sketch for \mathcal{F} . Then, $\ell_{\text{con}}(\mathcal{G}) \leq (1 + \delta)(1 + \zeta/2) \max(y, \ell_{\text{con}}(\mathcal{F}))$.

Proof: Assume that $\mathcal{F} = \{f_1, \dots, f_m\}$, and $\mathcal{G} = \{f_1, \dots, f_k\}$, where $k \leq m$. If $m = 1$ then $k = 1$, and by definition $\ell_{\text{con}}(\mathcal{F}) = \ell_{\text{con}}(\mathcal{G}) = 0$, and the result holds. If $m > 1$, we need to show that $(f_i)_{\leq y'}$, for $i = 1, \dots, k$, are connected, where $y' = (1 + \delta)(1 + \zeta/2)\ell$, and $\ell = \max(y, \ell_{\text{con}}(\mathcal{F}))$. By definition, $\mathcal{F}_{\leq \ell}$ is a connected set. Consider any distinct i and j with values between 1 and k . Then there is a sequence of distinct indices $i = i_1, i_2, \dots, i_s = j$ such that $(f_{i_r})_{\leq \ell} \cap (f_{i_{r+1}})_{\leq \ell} \neq \emptyset$, for $r = 1, \dots, s - 1$. Consider such an index, say i_r , such that $i_r > k$, i.e., $f_{i_r} \notin \mathcal{G}$. Since,

$$U = (f_{i_r})_{\leq \ell} \cap (f_{i_{r-1}})_{\leq \ell} \neq \emptyset \quad \text{and} \quad U' = (f_{i_r})_{\leq \ell} \cap (f_{i_{r+1}})_{\leq \ell} \neq \emptyset$$

one can choose two arbitrary points $\mathbf{u} \in U$ and $\mathbf{v} \in U'$. Now the entire segment $\mathbf{uv} \subseteq (f_{i_r})_{\leq (1+\zeta/2)\ell}$, by Lemma B.3. Since $(1 + \zeta/2)\ell \geq y$, it follows by the definition of a (δ, y) sketch, see Definition 2.10, that $\mathbf{uv} \subseteq (f_{i_r})_{\leq (1+\zeta/2)\ell} \subseteq \mathcal{G}_{\leq (1+\zeta/2)(1+\delta)\ell}$. By Observation B.4, the sets in the minimal cover of \mathbf{uv} by the sublevel sets $(f_i)_{\leq (1+\zeta/2)(1+\delta)\ell}$, are connected. It follows that $(f_{i_r})_{\leq (1+\zeta/2)\ell}$ can be replaced by a subset of $\left\{ (f_i)_{\leq (1+\zeta/2)(1+\delta)\ell} \mid i = 1, \dots, k \right\}$, and the property of neighbor intersections is still valid in the chain. We replace each occurrence of the set $(f_{i_r})_{\leq (1+\zeta/2)\ell}$ for $i_r > k$, by the corresponding chain. It is easy to see that the resulting chain connects up $(f_{i_1})_{\leq (1+\zeta/2)(1+\delta)\ell}$ and $(f_{i_s})_{\leq (1+\zeta/2)(1+\delta)\ell}$. Now, removing duplicate elements does not affect the neighbor intersection property of the resulting chain. ■

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

Lemma B.6. Let $\mathcal{G} \subseteq \mathcal{F}$ be sets of functions, where \mathcal{G} is a (δ, y_0) -sketch for \mathcal{F} for some $\delta \geq 0$ and $y_0 \geq 0$. Let \mathbf{q} be a point such that $\text{d}(\mathbf{q}, \mathcal{F}) \geq y_0$. Then we have that $\text{d}(\mathbf{q}, \mathcal{G}) \leq (1 + \delta)\text{d}(\mathbf{q}, \mathcal{F})$.

Proof: Let $\ell = \text{d}(\mathbf{q}, \mathcal{F})$ and let $f \in \mathcal{F}$ be a witness to this distance; that is, $\mathbf{q} \in f_{\leq \ell}$. As $\ell \geq y_0$, we have that $f_{\leq \ell} \subseteq \bigcup_{g \in \mathcal{G}} g_{\leq (1+\delta)\ell}$ by the sketch property (Definition 2.10). As such there is some function $g \in \mathcal{G}$, such that $\mathbf{q} \in g_{\leq (1+\delta)\ell}$. It follows that $\text{d}(\mathbf{q}, g) \leq (1 + \delta)\text{d}(\mathbf{q}, \mathcal{F})$. ■