

Shape Fitting on Point Sets with Probability Distributions

Maarten Löffler
Utrecht University
loffler@cs.uu.nl

Jeff Phillips
Duke University
jeffp@cs.duke.edu

October 25, 2018

Abstract

A typical computational geometry problem begins: Consider a set P of n points in \mathbb{R}^d . However, many applications today work with input that is not precisely known, for example when the data is sensed and has some known error model. What if we do not know the set P exactly, but rather we have a probability distribution μ_P governing the location of each point $p \in P$?

Consider a set of (non-fixed) points P , and let μ_P be the probability distribution of this set. We study several measures (e.g. the radius of the smallest enclosing ball, or the area of the smallest enclosing box) with respect to μ_P . The solutions to these problems do not, as in the traditional case, consist of a single answer, but rather a distribution of answers. We hence describe a data structure, called an ε -quantization, that can approximate such a distribution within ε in $O(1/\varepsilon)$ space. We also extend this data structure to answer higher dimensional queries of μ_P (e.g. the length and width of the smallest enclosing box in \mathbb{R}^2).

Rather than compute a new data structure for each measure we are interested in, we can also compute a single data structure that allows us to answer many questions at once. This data structure, an (ε, α) -kernel, is based on α -kernel coresets and can be used to create approximate ε -quantizations for geometric problems involving extent measures.

Thirdly, we introduce a data structure that can answer questions of the type ‘what is the probability that point q is in the smallest enclosing ball of P ?’ For a given distribution μ_P and summarizing shape (e.g. the smallest enclosing ball), we define an ε -shape inclusion probability function to be a function that assigns to a query point $q \in \mathbb{R}^d$ a value that is at most ε away from the probability that q is contained in this summarizing shape of P . This results in a probability description more directly linked to the space that the input points live in.

We provide simple and efficient randomized algorithms for computing all of these data structures, which are easy to implement and practical. We provide some experimental results to assert this. We also provide more involved deterministic algorithms for ε -quantizations for problems involving shapes with bounded VC-dimension that run in time polynomial in n and $1/\varepsilon$.

1 Introduction

The input for a typical computational geometry problem is a set P of n points in \mathbb{R}^2 , or more generally \mathbb{R}^d . Traditionally, such a set of points is assumed to be known exactly, and indeed, in the 1980s and 1990s such an assumption was often justified because much of the input data was hand-constructed for computer graphics or simulations. However, in many modern applications the input is sensed from the real world, and such data is inherently imprecise. Therefore, there is a growing need for methods that are able to deal with imprecision.

An early model to quantify imprecision in geometric data, motivated by finite precision of coordinates, is ε -*geometry*, introduced by Guibas *et al.* [10]. In this model, the input is given by a traditional point set P , where the imprecision is modeled by a single extra parameter ε . The true point set is not known, but it is certain that for each point in P there is a point in the disk of radius ε around it. This model has proven fruitful and is still often used due to its simplicity. To name a few, Guibas *et al.* [11] define *strongly convex* polygons: polygons that are guaranteed to stay convex, even when the vertices are perturbed by ε . Bandyopadhyay and Snoeyink [3] compute the set of all potential simplices in \mathbb{R}^2 and \mathbb{R}^3 that could belong to the Delaunay triangulation. Held and Mitchell [13] and Löffler and Snoeyink [15] study the problem of preprocessing a set of imprecise points under this model, so that when the true points are specified later some computation can be done faster.

A more involved model for imprecision can be obtained by not specifying a single ε for all the points, but allowing a different radius for each point, or even other shapes of imprecision regions. This allows for modeling imprecision that comes from different sources, independent imprecision in different dimensions of the input, etc. This extra freedom in modeling comes at the price of more involved algorithmic solutions, but still many results are available. Nagai and Tokura [19] compute the union and intersection of all possible convex hulls to obtain bounds on any possible solution, as does Ostrovsky-Berman and Joskowicz [20] in a setting allowing some dependence between points. Van Kreveld and Löffler [23] study the problem of computing the smallest and largest possible values of several geometric extent measures, such as the diameter or the radius of the smallest enclosing ball, where the points are restricted to lie in given regions in the plane. Kruger [14] extends some of these results to higher dimensions.

However, some applications dealing with sensed data provide more information about the imprecision than just a region, and a probability distribution governing the expected location of each point may be available. In robotic mapping [8] careful error models are used to govern the laser range finder data. In data mining [2] original data is often perturbed by a known model for privacy preserving purposes. In databases [6] large data sets may be summarized as probability distributions to store them more compactly. The atoms of a protein structure have probabilistic distributions as determined by NMR spectroscopy reconstruction algorithms [22], rotamers, or other variability. Similarly, probability distribution models are produced for GIS data, data from sensor networks, astrological data, and many other sources. In these cases, the above threshold error models could be adapted to this data by choosing an error distance beyond which the probability is below a certain threshold. However, the solutions produced under the threshold error models depend heavily on the boundary cases of the error model, while it is reasonable to expect the points are more likely to appear near the “center” of the regions. Working directly with probability distributions can provide more accurate answers to geometric questions about such sets of points.

This paper studies the computation of extent measures on uncertain point sets governed by probability distributions. Unsurprisingly, directly using the probability distribution error model creates harder algorithmic problems, and many questions may be impossible to answer exactly under this model. But since the data is imprecise to begin with, it is also reasonable to construct approximate answers. Our algorithms have approximation guarantees with respect to the original distributions, not an approximation of them. Instead of reinventing computational geometry for probability distributions, this paper reduces problems on data governed by probability distributions to discrete and well-studied computational geometry problems on precise point sets.

1.1 Problem Statement

Let $\mu_p : \mathbb{R}^d \rightarrow \mathbb{R}^+$ describe the probability distribution of a point p where $\int_{x \in \mathbb{R}^d} \mu_p(x) dx = 1$. Let $\mu_P : \mathbb{R}^d \times \mathbb{R}^d \times \dots \times \mathbb{R}^d \rightarrow \mathbb{R}^+$ describe the distribution of a point set P by the joint probability over each $p \in P$. For

simplicity we refer to the space $\mathbb{R}^d \times \dots \times \mathbb{R}^d$ as \mathbb{R}^{dn} when it is a product of n d -dimensional spaces. For this paper we will assume $\mu_P(q_1, q_2, \dots, q_n) = \prod_{i=1}^n \mu_{p_i}(q_i)$, so the distribution for each point is independent, although for our randomized algorithms this restriction can be easily circumvented.

Given a distribution μ_P we can ask questions traditionally asked of point sets that are given precisely, instead of as distributions (e.g. the diameter or the axis-aligned bounding box). In the presence of imprecision, the answer to such a question is not a single value or structure, but also a *distribution* of answers. The point of this paper is not just how to answer geometric questions about these distributions, but how to concisely represent them.

ε -Quantizations. Let $f : \mathbb{R}^{dn} \rightarrow \mathbb{R}$ be a single-valued function on a fixed point set, such as the radius of the minimum enclosing ball. For a query value v ,

$$f_{\mu_P}^{\leq}(v) = \int_{Q \in \mathbb{R}^{dn}} 1(f(Q) \leq v) \cdot \mu_P(Q) dQ,$$

where Q is taken over all size n point sets in \mathbb{R}^d and $1(\cdot)$ is the indicator function, is the probability that f will yield a value less than or equal to v , given the distribution μ_P . Then $f_{\mu_P}^{\leq}$ is the cumulative density function of the distribution of possible values that f can take. Ideally, we would return the function $f_{\mu_P}^{\leq}$ so we could quickly answer any query exactly, however, it is not clear how to compute closed forms for such functions for one specific value, let alone all values. Rather, we introduce a data structure, which we call an ε -quantization, to answer such queries approximately and efficiently. For an isotonic function $f_{\mu_P}^{\leq}$ and any value v , an ε -quantization, R , guarantees that $|R(v) - f_{\mu_P}^{\leq}(v)| \leq \varepsilon$. Furthermore, the size of an ε -quantization is always dependent only on ε , not on $|P|$ or μ_P .

Sometimes a statistic for a point set has multiple values, such as the width of the minimum enclosing axis-aligned rectangle along the x -axis and the y -axis. For a function $f : \mathbb{R}^{dn} \rightarrow \mathbb{R}^k$ let

$$f_{\mu_P}^{\preceq}(v_1, \dots, v_k) = \int_{Q \in \mathbb{R}^{dn}} 1(f(Q) \preceq (v_1, \dots, v_k)) \cdot \mu(Q) dQ,$$

where for a point $p \in \mathbb{R}^k$ the operation $p \preceq (v_1, \dots, v_k)$ determines whether $p_i \leq v_i$ for each i , where p_i is the i th coordinate of p . Note that $f_{\mu_P}^{\preceq}$ must be isotonic in the sense that for two points $p, q \in \mathbb{R}^k$ if $p \preceq q$ then $f_{\mu_P}^{\preceq}(p) \leq f_{\mu_P}^{\preceq}(q)$. A k -variate ε -quantization R for an isotonic function $f_{\mu_P}^{\preceq} : \mathbb{R}^k \rightarrow [0, 1]$ and for a query $v \in \mathbb{R}^k$ guarantees $|R(v) - f_{\mu_P}^{\preceq}(v)| \leq \varepsilon$. The size of a multivariate ε -quantization is dependent only on ε and k .

(ε, α) -Kernels. Rather than compute a new data structure for each measure we are interested in, we can also compute a single data structure that allows us to answer many questions at once. For an isotonic function $f_{\mu_P}^{\leq} : \mathbb{R}^+ \rightarrow [0, 1]$ an (ε, α) -quantization M guarantees that there exists a point x' such that (1) $|x - x'| \leq \alpha x$ and (2) $|M(x) - f_{\mu_P}^{\leq}(x')| \leq \varepsilon$. An (ε, α) -kernel is a data structure that can produce an (ε, α) -quantization for $f_{\mu_P}^{\leq}$ where f measures the width in any directions and whose size depends only on $\frac{1}{\varepsilon}$ and $\frac{1}{\alpha}$.

Shape Inclusion Probabilities. To summarize a point set $P \subset \mathbb{R}^d$, we often approximate it with a shape, such as the smallest enclosing ball. For k -variate ε -quantizations with large k , it can be hard to visualize the connection to the ambient d -dimensional space of the data points (i.e. for smallest enclosing ball we could use a $(d + 1)$ -variate ε -quantization to measure the d coordinates of the center point and the radius). Instead, for a summarizing shape we may wish to study a *shape inclusion probability function* $h_{\mu_P} : \mathbb{R}^d \rightarrow [0, 1]$ (or sip function) which describes the probability that a given point $x \in \mathbb{R}^d$ is included in the summarizing shape¹. Again, there does not seem to be a closed form for many of these functions. Rather we calculate an ε -sip function $\hat{h} : \mathbb{R}^d \rightarrow [0, 1]$ such that $\forall_{x \in \mathbb{R}^d} |h(x) - \hat{h}(x)| \leq \varepsilon$. The size of an ε -sip depends only on ε and the complexity of the summarizing shape.

¹For technical reasons, if there are (degenerately) multiple optimal summarizing shapes, we say each are equally likely to be the summarizing shape of the point set.

1.2 Our Results

We describe simple and practical randomized algorithms for computing ε -quantizations, ε -sip functions, and (ε, α) -kernels. Let $T_f(n)$ be the time it takes to calculate a summarizing shape of a set of n points $Q \subset \mathbb{R}^d$, which generates a statistic $f(Q)$. We can calculate an ε -quantization of $f_{\mu_P}^{\leq}$, with probability $1 - \delta$, in time $O(T_f(n) \frac{1}{\varepsilon^2} \log \frac{1}{\varepsilon \delta})$. For univariate ε -quantizations the size is $O(\frac{1}{\varepsilon})$, and for k -variate ε -quantizations the size is $O(k^2 \frac{1}{\varepsilon} \log^{2k} \frac{1}{\varepsilon})$. With probability $1 - \delta$, we can calculate an ε -sip function of size $O(\frac{1}{\varepsilon^2} \log \frac{1}{\varepsilon \delta})$ in time $O(T_f(n) \frac{1}{\varepsilon^2} \log \frac{1}{\varepsilon \delta})$. With probability $1 - \delta$, we can calculate an (ε, α) -kernel of size $O(\frac{1}{\alpha^{(d-1)/2} \varepsilon^2} \log \frac{1}{\varepsilon \delta})$ in time $O((n + \frac{1}{\alpha^{d-3/2}}) \frac{1}{\varepsilon^2} \log \frac{1}{\varepsilon \delta})$. All of these randomized algorithms are simple and practical, as demonstrated by some experimental results.

In addition, we provide deterministic algorithms for computing ε -quantizations of a specific class of functions. If \mathfrak{A} is a family of geometric shapes, such that $(\mathbb{R}^d, \mathfrak{A})$ has bounded VC-dimension, and $f : \mathbb{R}^{dn} \rightarrow \mathbb{R}^k$ is a function that describes some statistics on the smallest element from \mathfrak{A} that encloses the points (e.g. the radius of the smallest enclosing ball), then an ε -quantization for f can be computed in deterministic time $O(\text{poly}(n, \frac{1}{\varepsilon}))$, as described in Table 1.

This paper describes results for shape fitting problems for distributions of point sets in \mathbb{R}^d , in particular, we will use the smallest enclosing ball and the axis-aligned bounding box as running examples in the algorithm descriptions. We believe, though, that the concept of ε -quantizations should extend to many other problems with uncertain data. In fact, variations of our randomized algorithm should work for a more general array of problems.

2 Preliminaries: ε -Samples and α -Kernels

ε -Samples. For a set P (in our context a point set), let \mathfrak{A} be a set of subsets of P which for instance could be induced by containment in a shape from some family of geometric shapes. The pair (P, \mathfrak{A}) is called a *range space*. We say that Q is an ε -sample of (P, \mathfrak{A}) if

$$\forall R \in \mathfrak{A} \left| \frac{\phi(R \cap Q)}{\phi(Q)} - \frac{\phi(R \cap P)}{\phi(P)} \right| \leq \varepsilon,$$

where $|\cdot|$ takes the absolute value and $\phi(\cdot)$ returns the measure of a point set. In the discrete case $\phi(Q)$ returns the cardinality of Q . We say \mathfrak{A} *shatters* a set S if every subset of P is equal to $R \cap S$ for some $R \in \mathfrak{A}$. The cardinality of the largest discrete set $X \subseteq P$ that \mathfrak{A} can shatter is known as the *VC-dimension* of (P, \mathfrak{A}) .

When (P, \mathfrak{A}) has constant VC-dimension, we can create an ε -sample Q of (P, \mathfrak{A}) , with probability $1 - \delta$, by uniformly sampling $O(\nu \frac{1}{\varepsilon^2} \log \frac{\nu}{\delta \varepsilon})$ points from P [24]. There exist deterministic techniques to create ε -samples [16, 5] of size $O(\nu \frac{1}{\varepsilon^2} \log \frac{1}{\varepsilon})$ in time $O(\nu^3 \nu n (\frac{1}{\varepsilon^2} \log \frac{\nu}{\varepsilon})^\nu)$. When P is a point set in \mathbb{R}^d and the family of ranges \mathfrak{Q}_k is determined by inclusion of convex shapes whose sides have one of k predefined normal directions, such as the set of axis-aligned boxes, then an ε -sample for (P, \mathfrak{Q}_k) of size $O(\frac{k}{\varepsilon} \log^{2k} \frac{1}{\varepsilon})$ can be constructed in $O(\frac{n}{\varepsilon^3} \log^{6k} \frac{1}{\varepsilon})$ time [21].

When we have a distribution $\mu : \mathbb{R}^d \rightarrow \mathbb{R}^+$, such that $\int_{x \in \mathbb{R}^d} \mu(x) dx = 1$, we can think of this as the set P of points in \mathbb{R}^d , where the weight w of a point $p \in \mathbb{R}^d$ is $\mu(p)$. To simplify notation, we write (μ, \mathfrak{A}) as a range space where the ground set is this set $P = \mathbb{R}^d$ weighted by the distribution μ . Let it have VC-dimension ν . For distribution μ that is polygonally approximable [21] with a constant number of facets, we can construct an ε -sample of size $O(\frac{\nu}{\varepsilon^2} \log \frac{\nu}{\varepsilon})$ in time $O(\frac{\nu}{\varepsilon^2} \log^2 \frac{\nu}{\varepsilon})$. A longer primer on ε -samples is in Appendix A.

α -Kernels. Given a point set $P \in \mathbb{R}^d$ of size n and a direction $u \in \mathbb{S}^{d-1}$, let $P[u] = \arg \max_{p \in P} \langle p, u \rangle$, where $\langle \cdot, \cdot \rangle$ is the inner product operator. Let $\omega(P, u) = \langle P[u] - P[-u], u \rangle$ describe the width of P in direction u . We say that $K \subseteq P$ is an α -kernel of P if for all $u \in \mathbb{S}^{d-1}$

$$\omega(P, u) - \omega(K, u) \leq \alpha \cdot \omega(P, u).$$

α -kernels of size $O(\frac{1}{\alpha^{(d-1)/2}})$ can be calculated in time $O(n + \frac{1}{\alpha^{d-3/2}})$ [4]. Computing many extent related problems such as diameter and smallest enclosing ball on the α -kernel approximates the function on the original set [1].

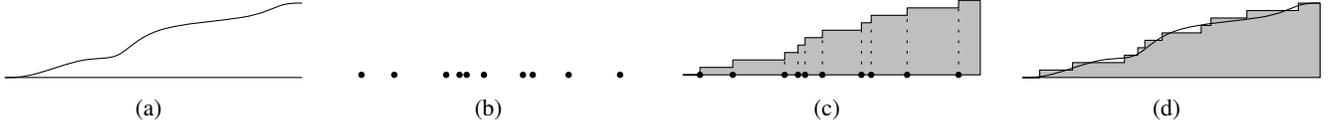


Figure 1: (a) The true form of the function. (b) The ε -quantization as a point set in \mathbb{R} . (c) The inferred curve in \mathbb{R}^2 . (d) Overlay of the two images.

3 Randomized Algorithm for ε -Quantizations

We start with a general algorithm (Algorithm 3.1) which will be made specific in several places in the paper. We assume we can draw a point from μ_p for each $p \in P$ in constant time; if the time depends on some other parameters, the runtimes can be easily adjusted.

Algorithm 3.1 Approximate μ_P with regard to a family of shapes \mathfrak{S} or function $f_{\mathfrak{S}}$

- 1: **for** $i = 1$ **to** $m = O(\frac{1}{\varepsilon^2} \log \frac{1}{\varepsilon\delta})$ **do**
 - 2: **for** $p_j \in P$ **do**
 - 3: Generate $q_j \in \mu_{p_j}$.
 - 4: Set $V_i = f_{\mathfrak{S}}(\{q_1, q_2, \dots, q_n\})$.
 - 5: Reduce or Simplify the set $V = \{V_i\}_{i=1}^m$.
-

Defining ε -quantizations. For an isotonic function $h : \mathbb{R} \rightarrow [0, 1]$, an ε -quantization, R , is a set of points where for any $t \in \mathbb{R}$, $|h(t) - R(t)| \leq \varepsilon$. We let $R(t) = \frac{1}{|R|} \sum_{r \in R} 1(r \leq t)$. Since h has range $[0, 1]$ and is isotonic, an ε -quantization requires only $O(1/\varepsilon)$ points. Figure 1 shows a illustration of how an ε -quantization approximates a smooth function. Because h is isotonic there exists a function $g : \mathbb{R} \rightarrow \mathbb{R}^+$ such that $h(t) = \int_{x=-\infty}^t g(x) dx$ where $\int_{x=-\infty}^{\infty} g(x) dx = 1$. Thus an ε -sample of (g, \mathbb{I}_+) is an ε -quantization of h , where \mathbb{I}_+ is all 1-sided intervals.

For an isotonic function $h : \mathbb{R}^k \rightarrow [0, 1]$ a k -variate ε -quantization, R , is a set of points in \mathbb{R}^k such that for any $p \in \mathbb{R}^k$, $|h(p) - R(p)| \leq \varepsilon$. For $p \in \mathbb{R}^k$ let $R(p) = \frac{1}{|R|} \sum_{q \in R} 1(q \preceq p)$. Because h is isotonic, there exists a function $g : \mathbb{R}^k \rightarrow \mathbb{R}^+$ such that $h(p) = \int_{x \preceq p} g(x) dx$ and $\int_{x \in \mathbb{R}^k} g(x) dx = 1$. Thus an ε -sample of (g, \mathbb{R}_+) is an ε -quantization of h , where \mathbb{R}_+ describes ranges $R_p \in \mathbb{R}_+$ defined by all q such that $q \preceq p$ for any p . See Figure 2 for an illustration of k -variate function h and a k -variate ε -quantization approximating it.

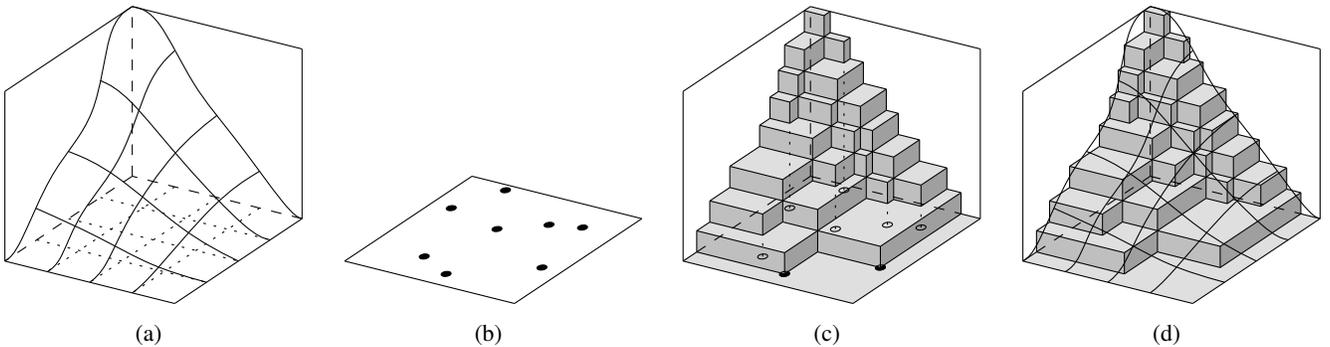


Figure 2: (a) The true form of the multivariate function. (b) The ε -quantisation as a point set in k -space. (c) The inferred surface in $k + 1$ -space. (d) Overlay of the two images.

Algorithm for ε -quantizations. For a function f on a point set P of size n , it takes $T_f(n)$ time to evaluate $f(P)$. We now construct $f_{\mu_P}^{\leq}$ by adapting Algorithm 3.1 as follows. First draw a sample point q_j from each μ_{p_j} for $p_j \in P$, then evaluate $V_i = f(\{q_1, \dots, q_n\})$. The fraction of trials of this process that produces a value less than v is the estimate of $f_{\mu_P}^{\leq}(v)$. Finally reduce the size of V by returning $\frac{2}{\varepsilon}$ evenly spaced points according to the sorted order.

Theorem 3.1. *For a distribution μ_P of n points, there exists a univariate ε -quantization of size $O(\frac{1}{\varepsilon})$ for $f_{\mu_P}^{\leq}$, and it can be constructed in $O(T_f(n)\frac{1}{\varepsilon^2} \log \frac{1}{\varepsilon\delta})$ time, with success probability $1 - \delta$, where $T_f(n)$ is the time it takes to compute $f(Q)$ for any point set Q of size n .*

Proof. Because $f_{\mu_P}^{\leq} : \mathbb{R} \rightarrow [0, 1]$ is an isotonic function, there exists another function $g : \mathbb{R} \rightarrow \mathbb{R}^+$ such that $f_{\mu_P}^{\leq}(t) = \int_{x=-\infty}^t g(x) dx$ where $\int_{\mathbb{R}} g(x) dx = 1$. And thus an ε -sample of (g, \mathbb{I}_+) is an ε -quantization of $f_{\mu_P}^{\leq}$.

By drawing a random sample q_i from each μ_{p_i} for $p_i \in P$, we are drawing a random point set Q from μ_P . Thus $f(Q)$ is a random sample from g . Hence, using the standard randomized construction for ε -samples, $O(\frac{1}{\varepsilon^2} \log \frac{1}{\varepsilon\delta})$ such samples will generate an $\frac{\varepsilon}{2}$ -sample for g , and hence an $\frac{\varepsilon}{2}$ -quantization for $f_{\mu_P}^{\leq}$, with probability $1 - \delta$.

Since in an $\frac{\varepsilon}{2}$ -quantization, every value is off from the true function by at most $\frac{\varepsilon}{2}$, then we can take an $\frac{\varepsilon}{2}$ -quantization of the step function and still have an ε -quantization of the true function. Thus, we can reduce this to an ε -quantization of size $O(\frac{1}{\varepsilon})$ by taking a subset of $\frac{2}{\varepsilon}$ points spaced evenly according to their sorted order. \square

We can construct k -variate ε -quantizations using the same basic procedure as in Algorithm 3.1. The output V_i of $f_{\mathfrak{S}}$ is k -variate and thus results in a k -dimensional point. As a result, the reduction of the final size of the point set requires more advanced procedures.

Theorem 3.2. *For a distribution μ_P of n points, there exists a k -variate ε -quantization of size $O(\frac{k^2}{\varepsilon} \log^{2k} \frac{1}{\varepsilon})$ for $f_{\mu_P}^{\leq}$, and it can be constructed in $O(T_f(n)\frac{k}{\varepsilon^2} \log \frac{k}{\varepsilon\delta} + k^2 \frac{1}{\varepsilon^5} \log^{6k} \frac{1}{\varepsilon} \log \frac{1}{\varepsilon\delta})$ time, with success probability $1 - \delta$, where $T_f(n)$ is the time it takes to compute $f(Q)$ for any point set Q of size n .*

Proof. In the k -variate case there exists a function $g : \mathbb{R}^k \rightarrow \mathbb{R}^+$ such that $f_{\mu_P}^{\leq}(v) = \int_{x \preceq v} g(x) dx$ where $\int_{\mathbb{R}^k} g(x) dx = 1$. Then a random point set Q from μ_P , evaluated as $f(Q)$, is still a random sample from the k -variate distribution described by g . Thus, with probability $1 - \delta$, a set of $O(\frac{k}{\varepsilon^2} \log \frac{1}{\varepsilon\delta})$ such samples is an ε -sample of (g, \mathbb{R}_+^k) , which has VC-dimension k , and the samples are also a k -variate ε -quantization of $f_{\mu_P}^{\leq}$.

We can then reduce the size of the ε -quantization to $O(\frac{k}{\varepsilon} \log^{2k} \frac{1}{\varepsilon})$ [21] (or to $O(\frac{k}{\varepsilon^2} \log \frac{1}{\varepsilon})$ [5]), since the VC-dimension is k and each data point requires $O(k)$ storage. \square

4 (ε, α) -Kernels

The above construction works for a fixed family of summarizing shapes. This section builds a single data structure, an (ε, α) -kernel, for a distribution μ_P in \mathbb{R}^d that can be used to construct (ε, α) -quantizations for several families of summarizing shapes. In particular, an (ε, α) -kernel of μ_P is a data structure such that in any query direction $u \in \mathbb{S}^{d-1}$ we can create an (ε, α) -quantization of $\omega(\cdot, u)$, the width in direction u . This data structure introduces a parameter α , which deals with geometric error, in addition to the error parameter ε , which deals with probability error.

We follow the randomized framework described above as follows. Let \mathfrak{K} be an (ε, α) -kernel consisting of $m = O(\frac{1}{\varepsilon^2} \log \frac{1}{\varepsilon\delta})$ α -kernels, where each α -kernel K_j approximates a point set Q_j drawn randomly from μ_P . Given \mathfrak{K} , we can then create an (ε, α) -quantization for the width of μ_P in any direction $u \in \mathbb{S}^{d-1}$. Specifically, let $M = \{\omega(K_j, u)\}_{j=1}^m$.

Lemma 4.1. *With probability $1 - \delta$, M is an (ε, α) -quantization of the width of μ_P in direction u .*

Proof. The width $\omega(Q_j, u)$ of a random point set Q_j drawn from μ_P is a random sample from the distribution over widths of μ_P in direction u . Thus, with probability $1 - \delta$, m such random samples would create an ε -quantization.

Using the width of the α -kernels K_j instead of Q_j induces an error on each random sample of at most $\alpha \cdot \omega(Q_j, u)$. Then for a query width w , say there are γm point sets Q_j that have width $\leq w$ and $\gamma' m$ α -kernels K_j with width $\leq w$. Note that $\gamma' > \gamma$. Let $\hat{w} = w - \alpha w$. For each point set Q_j that has width $> w$ but the corresponding α -kernel K_j has width $\leq w$, it follows that K_j has width $> \hat{w}$. Thus the number of α -kernels K_j that have width $\leq \hat{w}$ is $\leq \gamma m$, and thus there is a width w' between w and \hat{w} such that the number of α -kernels $\leq w'$ is exactly γm . \square

Theorem 4.1. *With probability $1 - \delta$, we can construct an (ε, α) -kernel for μ_P on n points in \mathbb{R}^d of size $O(\frac{1}{\alpha^{(d-1)/2}} \frac{1}{\varepsilon^2} \log \frac{1}{\varepsilon \delta})$ and in time $O((n + \frac{1}{\alpha^{d-3/2}}) \frac{1}{\varepsilon^2} \log \frac{1}{\varepsilon \delta})$.*

k -Dependent (ε, α) -Kernels. The definition of (ε, α) -quantizations can be extended to k -variate (ε, α) -quantizations M where (1) there exists a point $x' \in \mathbb{R}^k$ such that for all integers $i \in [1, k]$ $|x^{(i)} - (x')^{(i)}| \leq \alpha x^{(i)}$ and (2) $|M(x) - f_{\mu_P}^{\preceq}(x')| \leq \varepsilon$. Let $x^{(i)}$ represent the i th coordinate of a point $x \in \mathbb{R}^k$.

(ε, α) -kernels can be generalized to approximate other functions $f : \mathbb{R}^{dn} \rightarrow \mathbb{R}^k$, specified as follows. We say a point $p' \in \mathbb{R}^k$ is a *relative θ -approximation* of $p \in \mathbb{R}^k$ if for each coordinate i we have $p^{(i)} - p'^{(i)} \leq \theta p^{(i)}$. For functions f and θ where $f(K)$ is a relative $\theta(\alpha)$ -approximation of $f(Q)$ when K is an α -kernel of Q , we say that f is *relative $\theta(\alpha)$ -approximable*.

By setting $m = O(\frac{k}{\varepsilon^2} \log \frac{k}{\varepsilon \delta})$ in the above algorithm, with probability $1 - \delta$, we can build a k -dependent (ε, α) -kernel data structure \mathbb{K} with the following properties. It has size $O(\frac{1}{\alpha^{(d-1)/2}} \frac{k}{\varepsilon^2} \log \frac{k}{\varepsilon \delta})$ and can be built in time $O((n + \frac{1}{\alpha^{d-3/2}}) \frac{k}{\varepsilon^2} \log \frac{k}{\varepsilon \delta})$. To create a k -variate (ε, α) -quantization for a function f , create a k -dimensional point $p_j = f(K_j)$ for each α -kernel K_j in \mathbb{K} . The set M of m k -dimensional points forms the k -variate (ε, α) -quantization.

Theorem 4.2. *Let f be a relative $\theta(\alpha)$ -approximable function that takes $T_f(n)$ time to evaluate on a set of n points in \mathbb{R}^d . From a k -dependent (ε, α) -kernel \mathbb{K} with m α -kernels, with probability $1 - \delta$, we can create a k -variate $(\varepsilon, \theta(\alpha))$ -quantization of f , of size $O(\frac{1}{\varepsilon} \log^{2k} \frac{1}{\varepsilon})$ in time $O(T_f(\frac{1}{\alpha^{(d-1)/2}})m)$.*

Proof. Each evaluation of f on a point set Q_j drawn from μ_P is a random sample from the distribution over f on point sets drawn from μ_P and hence these values on all m sampled point sets would be an ε -quantization of $f_{\mu_P}^{\preceq}$.

For a query point $w \in \mathbb{R}^k$, let γm point sets produce a value $w_j = f(Q_j)$ such that $w_j \preceq w$, and let $\gamma' m$ point sets produce a value $w'_j = f(K_j)$ such that $w'_j \preceq w$. Note that $\gamma' > \gamma$. Because f is relative $\theta(\alpha)$ -approximable, for each point set Q_j such that $w_j \preceq w$, but $w_j \not\preceq \hat{w}$, then $w'_j \not\preceq \hat{w}$, where $\hat{w} = w - \theta(\alpha)w$. (More specifically, for each coordinate $w^{(i)}$ of w , $\hat{w}^{(i)} = w^{(i)} - \theta(\alpha)w^{(i)}$.) Thus, the number of point sets such that $f(K_j) \preceq \hat{w}$ is $\leq \gamma m$, and hence there is a point w' between w and \hat{w} such that the fraction of sampled point sets such that $f(K_j) \preceq w'$ is exactly γ , and hence is within ε of the true fraction of point sets sampled from μ_P with probability $1 - \delta$. \square

To name a new examples, the width and diameter are relative α -approximable functions, thus the results apply directly with $k = 1$. The radius of the minimum enclosing ball is relative 2α -approximable with $k = 1$. The d directional widths of the minimum perimeter or minimum volume axis-aligned rectangle is relative α -approximable with $k = d$.

4.1 Experiments with (ε, α) -Kernels and ε -Quantizations

We implemented these randomized algorithms for (ε, α) -kernels and ε -quantizations for diameter (diam), width in a fixed direction (dwid), and radius of the smallest enclosing ℓ_2 ball (seb₂). We used existing code from Hai Yu [25] for α -kernels and Bernd Gärtner [9] for seb₂. For the input set μ_P we generated 5000 points $P \subset \mathbb{R}^3$ on the surface of a cylinder piece with radius 1 and axis length 10. Each point $p \in P$ represented the center of a Gaussian with standard deviation 3. We set $\varepsilon = .2$ and generated α -kernels of size at most 40 (the existing code did not allow the use to specify a parameter α , only the maximum size). We generated a total of $m = 40$ point sets from μ_P . The (ε, α) -kernel has a total of 1338 points. We calculated ε -quantizations and (ε, α) -quantizations for diam, dwid, and seb₂, each of size 10; see Figure 3.

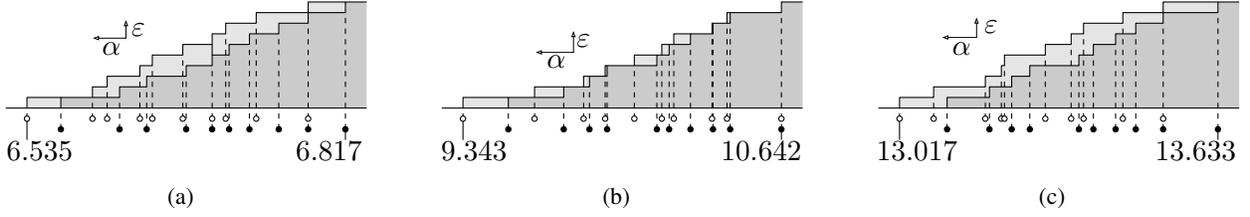


Figure 3: (ε, α) -quantization (white circles) and ε -quantization (black circles) for (a) seb_2 , (b) dwid , and (c) diam .

5 Shape Inclusion Probabilities

We can also use a variation of Algorithm 3.1 to construct ε -shape inclusion probability functions. For a point set $Q \subset \mathbb{R}^d$, let the summarizing shape $S_Q = \mathfrak{S}(Q)$ be from some geometric family \mathfrak{S} so $(\mathbb{R}^d, \mathfrak{S})$ has bounded VC-dimension ν . We randomly sample point sets Q_j from μ_P and then find the summarizing shape S_{Q_j} (e.g. minimum enclosing ball) of Q_j . Let this set of shapes be $S^{(\mu_P)}$. If there are multiple shapes from \mathfrak{S} which are equally optimal (as can happen degenerately with, for example, minimum width slabs), choose one of these shapes at random. For a set of shapes $S' \subset \mathfrak{S}$, let $S'_p \subset S'$ be the subset of shapes that contain $p \in \mathbb{R}^d$. We store $S^{(\mu_P)}$ and evaluate a query point $p \in \mathbb{R}^d$ by counting what fraction of the shapes the point is contained in, specifically returning $|S'_p|/|S^{(\mu_P)}|$ in $O(\nu|S^{(\mu_P)}|)$ time. In some cases, this evaluation can be sped up with point location data structures.

Theorem 5.1. *For a distribution μ_P of n points and a family of summarizing shapes $(\mathbb{R}^d, \mathfrak{S})$ with bounded VC-dimension ν , with probability $1 - \delta$ we can construct an ε -sip function of size $O(2^{\nu+1} \frac{\nu^2}{\varepsilon^2} \log \frac{1}{\varepsilon\delta})$ and in time $O(T_{\mathfrak{S}}(n) \frac{1}{\varepsilon^2} \log \frac{1}{\varepsilon\delta})$, where $T_{\mathfrak{S}}(n)$ is the time it takes to determine the summarizing shape of any point set $Q \subset \mathbb{R}^d$ of size n .*

Proof. If $(\mathbb{R}^d, \mathfrak{S})$ has VC-dimension ν , then the dual range space (\mathfrak{S}, P^*) has VC-dimension $\nu' \leq 2^{\nu+1}$, where P^* is all subsets $\mathfrak{S}_p \subseteq \mathfrak{S}$, for any $p \in \mathbb{R}^d$, such that $\mathfrak{S}_p = \{S \in \mathfrak{S} \mid p \in S\}$. Using the above algorithm, sample $m = O(\frac{\nu'}{\varepsilon^2} \log \frac{\nu'}{\varepsilon\delta})$ point sets Q from μ_P and generate the m summarizing shapes S_{Q_j} . Each shape is a random sample from \mathfrak{S} according to μ_P , and thus $S^{(\mu_P)}$ is an ε -sample of (\mathfrak{S}, P^*) .

Let $w_{\mu_P}(S)$, for $S \in \mathfrak{S}$, be the probability that S is the summarizing shape of a point set Q drawn randomly from μ_P . Let $W_{\mu_P}(\mathfrak{S}') = \int_{S \in \mathfrak{S}'} w_{\mu_P}(S)$, where $\mathfrak{S}' \subseteq P^*$, be the probability that some shape from the subset \mathfrak{S}' is the summarizing shape of Q drawn from μ_P .

We approximate the sip function at $p \in \mathbb{R}^d$ by returning the fraction $|S'_p|/m$. The true answer to the sip function at $p \in \mathbb{R}^d$ is $W_{\mu_P}(\mathfrak{S}_p)$. Since $S^{(\mu_P)}$ is an ε -sample of (\mathfrak{S}, P^*) , then with probability $1 - \delta$

$$\left| \frac{|S'_p|}{m} - \frac{W_{\mu_P}(\mathfrak{S}_p)}{1} \right| = \left| \frac{|S'_p|}{|S^{(\mu_P)}|} - \frac{W_{\mu_P}(\mathfrak{S}_p)}{W_{\mu_P}(P^*)} \right| \leq \varepsilon.$$

Since for the family of summarizing shapes \mathfrak{S} the range space $(\mathbb{R}^d, \mathfrak{S})$ has VC-dimension ν , each can be stored using that much space. \square

The size can then be reduced to $O(2^{\nu+1} \frac{\nu^2}{\varepsilon^2} \log \frac{1}{\varepsilon})$ in time $O((2^{\nu+1})^{3 \cdot 2^{\nu+1} + 1} (\frac{\nu}{\varepsilon} \log \frac{1}{\varepsilon})^{2^{\nu+1} + 1})$ using deterministic techniques.

Representing ε -sip functions by Isolines. Shape inclusion probability functions are density functions. One convenient way of visually representing a density function in \mathbb{R}^2 is by drawing the isolines. A γ -isoline is a closed curve such that on the inside the density function is $> \gamma$ and on the outside is $< \gamma$.

In each part of Figure 4 a set of 5 circles correspond to points with a probability distribution. For part (a) and (c), the probability distribution μ_P is uniform over those circles, in part (b) and (d) it is drawn from a multivariate Gaussian

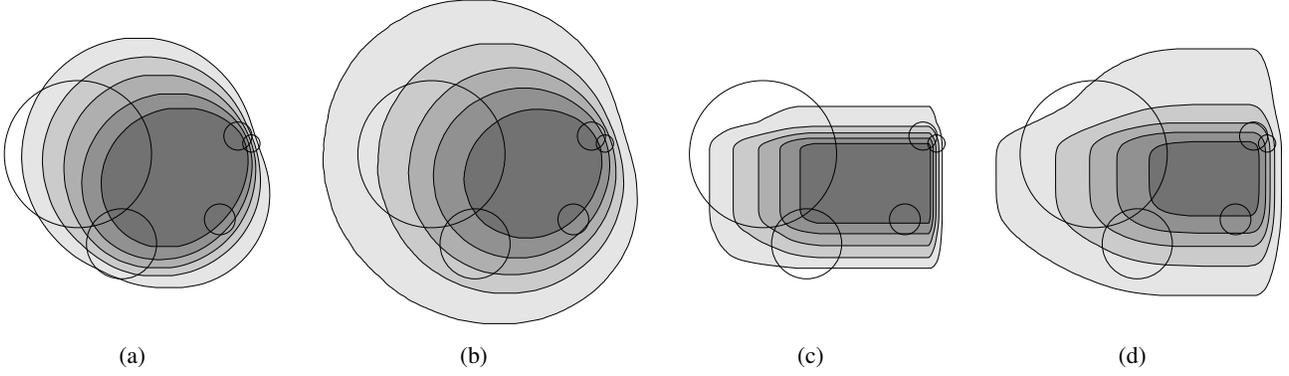


Figure 4: (a) The shape inclusion probability for the smallest enclosing ball, for points uniformly distributed inside the circles. (b) The same, but for normally distributed points around the circle centers, with standard deviations given by the radii. (c) The shape inclusion probability for the smallest enclosing axis-aligned rectangle, for points uniformly distributed inside the circles. (d) The same, but for normally distributed points.

distribution with standard deviation as the radius. We generate ε -sip functions for smallest enclosing ball in Figure 4(a,b) and for smallest axis-aligned bounding box in Figure 4(c,d).

In all figures we draw approximations of $\{.9, .7, .5, .3, .1\}$ -isolines. These drawings are generated by randomly selecting $m = 5000$ (a,b) or $m = 25000$ (c,d) shapes, counting the number of inclusions at different points in the plane and interpolating to get the isolines. The innermost and darkest region has probability $> 90\%$, the next one probability $> 70\%$, etc., the outermost region has probability $< 10\%$.

When μ_P describes the distribution for n points and n is large, then isolines are generally connected for convex summarizing shapes. In fact, in $O(n)$ time we can create a point which is contained in the convex hull of a point set sampled from μ_P with high probability. Specifics are discussed in Appendix B.

6 Deterministic Constructions of ε -Quantizations

In this section we consider functions f which describe the size of some summarizing shape from the family \mathbb{A} such that $(\mathbb{R}^d, \mathbb{A})$ has constant VC-dimension. In particular, given a point set $Q \subset \mathbb{R}^d$, let $\mathbb{A}(Q) \subset \mathbb{R}^d$ (e.g. smallest enclosing ball) be the summarizing shape for Q , and let $f(Q)$ be a statistic of $\mathbb{A}(Q)$ (e.g. radius of the smallest enclosing ball). The overall strategy will be to deterministically approximate each μ_{p_i} with a point set Q_{p_i} , although not with respect to the range space (μ_{p_i}, \mathbb{A}) , but with a more complicated range space described below. Let $Q_P = \{Q_{p_i}\}_i$ describe this set of point sets. Then let the function $f(Q_P, r)$ describe the fraction of point sets $Q' = (q_1 \in Q_{p_1}, q_2 \in Q_{p_2}, \dots, q_n \in Q_{p_n})$ for $\{Q_{p_1}, \dots, Q_{p_n}\} = Q_P$ such that $f(Q') \leq r$. We show that we can generate a set of point sets Q_P such that $f(Q_P, r)$ is a good approximation of $f_{\mu_P}^{\leq}(r)$. And we show how to efficiently evaluate $f(Q_P, r)$.

6.1 Approximating μ_p

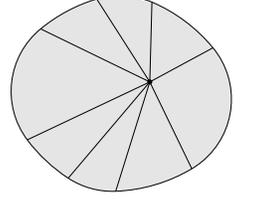
In this section we restrict that μ_P is either defined by a polygonal surface S with b facets or is polygonal approximable, it can be approximated by a finite polygonal surface S with b facets, for some constant b , as in [21].

It might seem that we can just create an ε -sample of (μ_{p_i}, \mathbb{A}) for each μ_{p_i} , but we need to consider a more complicated family $\mathbb{A}_{f,n}$. Given a family of shapes \mathbb{A} and a function f which computes a value determined by a summarizing shape $A \in \mathbb{A}$ for a set of n points, then $\mathbb{A}_{f,n}$ is a family of shapes where each is defined by a set of $n - 1$ points $T \subset \mathbb{R}^d$ and a value w . Specifically, $\mathbb{A}_{f,n}(T, w)$ is the set of points $\{p \in \mathbb{R}^d \mid f(T \cup p) \leq w\}$.

In certain cases, such as the volume of the axis-aligned bounding box, $(\mu_{p_i}, \mathbb{A}_{f,n})$ has constant VC-dimension. Shapes from $\mathbb{A}_{f,n}$ are determined by the placement of $2d$ points, the most extreme in each axis direction, thus its

shatter dimension is $\sigma_f = 2d$. Hence an ε -sample for $(\mu_{p_i}, \mathbb{A}_{f,n})$ of size $O(\frac{1}{\varepsilon^2} \log \frac{1}{\varepsilon})$ can be calculated in time $O(\frac{1}{\varepsilon^2} \log^2 \frac{1}{\varepsilon})$ for each μ_{p_i} .

In other cases, such as the radius of smallest enclosing ℓ_2 disks, $\mathbb{A}_{f,n}$ defines regions which have $O(n)$ $(d-1)$ -dimensional faces on its boundary and thus $(\mu_{p_i}, \mathbb{A}_{f,n})$ has VC-dimension n . Naive techniques would take time exponential in n to deterministically create an ε -sample, but we can do better by decomposing a shape $A \in \mathbb{A}_{f,n}$ into $O(n)$ disjoint simpler shapes. In the case of disks, $\mathbb{A}_{f,n}$ has its boundary defined by at most $2n$ circular arcs of two different radii. We can choose a point in the convex hull of T and draw lines to each intersection of circular arcs, see the figure on right. The intersections of the disc defining each boundary piece and the halfspaces for the drawn lines at its endpoints describes a *wedge* from a family $\mathbb{W}_{f,n}$. The range space $(\mathbb{R}^d, \mathbb{W}_{f,n})$ has VC-dimension at most 9 because shapes from $\mathbb{W}_{f,n}$ are formed by the intersection of three shapes from families that would each have VC-dimension 3 in a range space on the same ground set. Thus, an $\frac{\varepsilon}{2n}$ -sample of $(\mu_{p_i}, \mathbb{W}_{f,n})$ is an ε -sample of $(\mu_{p_i}, \mathbb{A}_{f,n})$. So for radius of the smallest enclosing ℓ_2 balls we can create an ε -sample of $(\mu_{p_i}, \mathbb{A}_{f,n})$ of size $O(n^2 \frac{1}{\varepsilon^2} \log \frac{n}{\varepsilon})$ in time $O(n^2 \frac{1}{\varepsilon^2} \log^2 \frac{n}{\varepsilon})$ for each μ_{p_i} .



We generalize both of these cases to other shapes and in higher dimensions in Appendix C. There are also illustrations of various shapes from $\mathbb{A}_{f,n}$.

Lemma 6.1. *When each μ_{p_i} is approximated with an ε' -sample Q_{p_i} of $(\mu_{p_i}, \mathbb{A}_{f,n})$, then for any r*

$$|\Pr[f_{\mu_P}(P) \leq r] - f(\{Q_{p_1}, Q_{p_2}, \dots, Q_{p_n}\}, r)| \leq \varepsilon' n.$$

Proof. When P is drawn from a distribution μ_P , then we can write that probability that $f_{\mu_P}(P) \leq r$ as follows.

$$\Pr[f_{\mu_P}(P) \leq r] = \int_{q_1} \mu_{p_1}(q_1) \int_{q_2} \mu_{p_2}(q_2) \dots \int_{q_n} \mu_{p_n}(q_n) \mathbf{1}(f(\{q_1, q_2, \dots, q_n\}) \leq r) dq_n dq_{n-1} \dots dq_1$$

Consider the inner most integral

$$\int_{q_n} \mu_{p_n}(q_n) \mathbf{1}(f(\{q_1, q_2, \dots, q_n\}) \leq r) dq_n,$$

where $\{q_1, q_2, \dots, q_{n-1}\}$ are fixed. The indicator function is true when for q_n $f(\{q_1, q_2, \dots, q_{n-1}, q_n\}) \leq r$ and hence q_n is contained in a shape from $\mathbb{A}_{f,n}(\{q_1, q_2, \dots, q_{n-1}\}, r)$. Thus if we have an ε' -sample Q_{p_n} for $(\mu_{p_n}, \mathbb{A}_{f,n})$, then we can guarantee that

$$\int_{q_n} \mu_{p_n}(q_n) \mathbf{1}(f(\{q_1, q_2, \dots, q_n\}) \leq r) dq_n \leq \frac{1}{|Q_{p_n}|} \sum_{q_n \in Q_{p_n}} \mathbf{1}(f(\{q_1, q_2, \dots, q_{n-1}, q_n\}) \leq r) + \varepsilon'.$$

We can then move the ε' to the outside, and we can change the order of the integrals to write:

$$\Pr[f_{\mu_P}(P) \leq r] \leq \frac{1}{|Q_{p_n}|} \sum_{q_n \in Q_{p_n}} \int_{q_1} \mu_{p_1}(q_1) \int_{q_2} \mu_{p_2}(q_2) \dots \int_{q_{n-1}} \mu_{p_{n-1}}(q_{n-1}) \mathbf{1}(f(\{q_1, q_2, \dots, q_n\}) \leq r) dq_{n-1} dq_{n-2} \dots dq_1 + \varepsilon'.$$

Repeating this procedure n times we get:

$$\begin{aligned} \Pr[f_{\mu_P}(P) \leq r] &\leq \left(\prod_{i=1}^n \frac{1}{|Q_{p_i}|} \right) \sum_{i=1}^n \sum_{q_i \in Q_{p_i}} \mathbf{1}(f(\{q_1, q_2, \dots, q_n\}) \leq r) + \varepsilon' n. \\ &= f(Q_P, r) + \varepsilon' n. \end{aligned}$$

Using the same technique we can achieve a symmetric lower bound for $\Pr[f_{\mu_P}(P) \leq r]$. □

By setting $\varepsilon' = \varepsilon/n$ we can achieve an additive ε -approximation by using an ε' -sample for each $(\mu_{p_i}, \mathbb{A}_{f,n})$.

6.2 Evaluating $f(Q_P, r)$.

Evaluating $f(Q_P, r)$ in time polynomial in n and $|Q_{p_i}|$, for any i , is not completely trivial since there are $n^{|Q_{p_i}|}$ possible sets in Q_P . Let a *good set* be a set of n points, G , such that for each Q_i there exists a point $g_i \in G$ such that $g_i \in Q_i$. For each good set G there exists a unique basis of at most σ_f points² which define the summarizing shape of G (remember the shatter dimension of \mathbb{A} is σ_f and $\sigma_f < \nu_f$, the VC-dimension). Define a *valid basis* to be a set of at most σ_f points in Q_P such that each point is from a different Q_i and if any point is removed the summarizing shape changes. Each valid basis forms a basis for several good sets.

We now construct R , an ε -quantization of $f_{\mu_P}^{\leq}$. This approximation is created by calculating the summarizing shape for all good sets. Even though there are an exponential number of good sets, there are only a polynomial number of valid bases. Thus for each valid basis, we count the number of good sets it represents. And we let each valid basis contribute to the ε -quantization; its position is determined by its value in f and its weight by the number of good sets it represents. We initially store the ε -quantization as a sorted list of tuples (r, η) where $r = f(\{q_1, q_2, \dots, q_{\sigma_f}\})$ for some valid basis $\{q_1, q_2, \dots, q_{\sigma_f}\}$, and η is the fraction of the good sets which are represented by this valid basis. The details are outlined in Algorithm 6.1.

Algorithm 6.1 Construct ε -Quantization from Q_P

```

1: for all valid bases  $q_1, q_2, \dots, q_{\sigma_f} \in Q_P$  do
2:   for  $i = 1$  to  $n$  do
3:     if  $q_1 \in Q_i$  or  $q_2 \in Q_i$  or  $\dots$  or  $q_{\sigma_f} \in Q_i$  then
4:       Set  $w_i = \frac{1}{|Q_i|}$ .
5:     else
6:       Set  $w_i = \frac{1}{|Q_i|} \sum_{q_j \in Q_i} 1(q_j \in \mathbb{A}(\{q_1, q_2, \dots, q_{\sigma_f}\}))$ 
7:   Insert  $(f(q_1, q_2, \dots, q_{\sigma_f}), \prod_i w_i)$  into  $R$ .
```

We now summarize the full deterministic algorithm. For each $(\mu_{p_i}, \mathbb{A}_{f,n})$ we create an $\frac{\varepsilon}{n}$ -sample Q_{p_i} of size $\alpha_f(n, \varepsilon)$. This makes the set Q_P have $\eta = \sum_{i=1}^n |Q_{p_i}| = n\alpha_f(n, \varepsilon)$ points in its sets. We examine $O(\eta^{\sigma_f})$ valid bases. For each valid basis we can evaluate $f(G)$ and compute w_i in $RS_f(n, \varepsilon)$ time using a range searching data structure, after preprocessing or with a naive search. Thus the deterministic running time for constructing an ε -quantization is $O(\eta^{\nu_f} RS_f(n, \varepsilon))$ which is presented for various summarizing shapes in Table 3. For instance, for volume of the axis-aligned bounding box this takes $O(n^{6d}/\varepsilon^{4d} \log^{3d} \frac{n}{\varepsilon})$ time and for radius of the smallest enclosing disks this takes $O(n^{16.5}/\varepsilon^7 \log^{3.5} \frac{n}{\varepsilon})$. The total construction time for the ε -quantizations is the sum of this time and the time to construct n (ε/n) -samples of $(\mathbb{R}^d, \mathbb{A}_{f,n})$; for both smallest enclosing disks and for axis-aligned bounding boxes it is the former.

A univariate ε -quantization can be reduced to size $O(\frac{1}{\varepsilon})$. Furthermore, we can create k -variate ε -quantizations using the same procedure (such as the width in the k dimensions of an axis-aligned bounding box). The condition in Lemma 6.1 where $f_{\mu_P}(P) \leq r$ can be replaced with a k -variate condition $f_{\mu_P}(P) \preceq r$ for $r \in \mathbb{R}^k$. Thus the same argument applies when we define $f : \mathbb{R}^{dn} \rightarrow \mathbb{R}^k$, and we can create k -variate ε -quantizations of size $k^2 \frac{1}{\varepsilon} \log^{O(k)} \frac{k}{\varepsilon}$ in the same deterministic times as long as $\nu_f = O(k)$.

Theorem 6.1. *For any range space (μ_P, \mathbb{A}_f) for a distribution μ_P of n points, with VC-dimension ν_f , where each $(\mu_{p_i}, \mathbb{A}_{f,n})$ has an $\frac{\varepsilon}{n}$ -sample of size $\alpha_f(n, \varepsilon)$, and where, after preprocessing m points and with near-linear space and time, we can count the number of points in a shape from \mathbb{A}_f in $RS(m, \mathbb{A}_f)$ time, we construct a k -variate ε -quantization of $f_{\mu_P}^{\leq}$ of size $k^2 \frac{1}{\varepsilon} \log^{O(k)} \frac{k}{\varepsilon}$ in $O((n\alpha_f(n, \varepsilon))^{\nu_f} \cdot RS((n\alpha_f(n, \varepsilon))^{\nu_f}, \mathbb{A}_f))$ time.*

²This uniqueness requires careful construction of the ε -samples Q_i , as described in Appendix A.1.

7 Acknowledgements

We would like to thank Pankaj K. Agarwal for many helpful discussions and Sariel Har-Peled for suggesting the use of wedges.

References

- [1] Pankaj K. Agarwal, Sariel Har-Peled, and Kasturi R. Varadarajan. Approximating extent measure of points. *Journal of ACM*, 51(4):2004, 2004.
- [2] Rakesh Agarwal and Ramakrishnan Srikant. Privacy-preserving data mining. *ACM SIGMOD Record*, 29:439–450, 2000.
- [3] Deepak Bandyopadhyay and Jack Snoeyink. Almost-Delaunay simplices: Nearest neighbor relations for imprecise points. In *ACM-SIAM Symp on Discrete Algorithms*, pages 403–412, 2004.
- [4] Timothy Chan. Faster core-set constructions and data-stream algorithms in fixed dimensions. *Computational Geometry: Theory and Applications*, 35:20–35, 2006.
- [5] Bernard Chazelle and Jiri Matousek. On linear-time deterministic algorithms for optimization problems in fixed dimensions. *Journal of Algorithms*, 21:579–597, 1996.
- [6] Reynold Cheng, Dmitri V. Kalashnikov, and Sunil Prabhakar. Evaluating probabilistic queries over imprecise data. In *Proceedings 2003 ACM SIGMOD International Conference on Management of Data*, 2003.
- [7] Kenneth L. Clarkson, David Eppstein, Gary L. Miller, Carl Sturtivant, and Shang-Hua Teng. Approximating center points with iterative Radon points. *International Journal of Computational Geometry and Applications*, 6:357–377, 1996.
- [8] Austin Eliazar and Ronald Parr. Dp-slam 2.0. In *Proceedings 2004 IEEE International Conference on Robotics and Automation*, 2004.
- [9] Bernd Gärtner. Fast and robust smallest enclosing balls. In *Proceedings 7th Annual European Symposium on Algorithms*, volume LNCS 1643, pages 325–338, 1999.
- [10] Leonidas J. Guibas, D. Salesin, and J. Stolfi. Epsilon geometry: building robust algorithms from imprecise computations. In *Proc. 5th Annu. ACM Sympos. Comput. Geom.*, pages 208–217, 1989.
- [11] Leonidas J. Guibas, D. Salesin, and J. Stolfi. Constructing strongly convex approximate hulls with inaccurate primitives. *Algorithmica*, 9:534–560, 1993.
- [12] Sariel Har-Peled. *Approximation Algorithm in Geometry*. <http://valis.cs.uiuc.edu/~sariel/teach/notes/aprx/>, 2008.
- [13] Martin Held and Joseph S. B. Mitchell. Triangulating input-constrained planar point sets. *Information Processing Letters*, page to appear, 2008.
- [14] Heinrich Kruger. Basic measures for imprecise point sets in \mathbb{R}^d . Master’s thesis, Utrecht University, 2008.
- [15] Maarten Löffler and Jack Snoeyink. Delaunay triangulations of imprecise points in linear time after preprocessing. In *Proc. 24th Symposium on Computational Geometry*, pages 298–304, 2008.
- [16] Jiri Matousek. Approximations and optimal geometric divide-and-conquer. In *Proceedings of the 23rd Annual ACM Symposium on Theory of Computing*, pages 505–511, 1991.

- [17] Jiri Matousek. *Geometric Discrepancy; An Illustrated Guide*, volume 18 of *Algorithms and Combinatorics*. Springer, 1999.
- [18] Jiri Matousek, Emo Welzl, and Lorenz Wernisch. Discrepancy and approximations for bounded vc-dimension. *Combinatorica*, 13(4):455–466, 1993.
- [19] T. Nagai and N. Tokura. Tight error bounds of geometric problems on convex objects with imprecise coordinates. In *Jap. Conf. on Discrete and Comput. Geom.*, LNCS 2098, pages 252–263, 2000.
- [20] Y. Ostrovsky-Berman and L. Jostkiewicz. Uncertainty envelopes. In *Abstracts 21st European Workshop on Comput. Geom.*, pages 175–178, 2005.
- [21] Jeff M. Phillips. Algorithms for ε -approximations of terrains. In *Proceedings 35th International Colloquium on Automata, Languages, and Programming*, 2008. arXiv 0801.2793.
- [22] Shobha Potluri, Anthony K. Yan, James J. Chou, Bruce R. Donald, and Chris Baily-Kellogg. Structure determination of symmetric homo-oligomers by complete search of symmetry configuration space, using nmr restraints and van der Waals packing. *Proteins*, 65:203–219, 2006.
- [23] Marc van Kreveld and Maarten Löffler. Largest bounding box, smallest diameter, and related problems on imprecise points. In *Proc. 10th Workshop on Algorithms and Data Structures*, LNCS 4619, pages 447–458, 2007.
- [24] Vladimir Vapnik and Alexey Chervonenkis. On the uniform convergence of relative frequencies of events to their probabilities. *Theory of Probability and its Applications*, 16:264–280, 1971.
- [25] Hai Yu, Pankaj K. Agarwal, Raghunath Poreddy, and Kasturi R. Varadarajan. Practical methods for shape fitting and kinetic data structures using coresets. In *Proceedings 20th Annual Symposium on Computational Geometry*, 2004.

A Primer on ε -Samples

We recall from Section 2 that for a range space (P, \mathfrak{A}) an ε -sample $Q \subseteq P$ guarantees

$$\forall R \in \mathfrak{A} \left| \frac{\phi(R \cap Q)}{\phi(Q)} - \frac{\phi(R \cap P)}{\phi(P)} \right| \leq \varepsilon,$$

where $|\cdot|$ takes the absolute value and $\phi(\cdot)$ returns the measure of a point set. In the discrete case $\phi(Q)$ returns the cardinality of Q .

When $P \subset \mathbb{R}^d$ we describe a few common examples of \mathfrak{A} . Let \mathfrak{B} describe all subsets of P determined by containment in some ball. Let \mathfrak{R}_d describe all subsets of P defined by containment in some d -dimensional axis-aligned box. Let \mathfrak{H} describe all subsets of P defined by containment in some halfspace. Throughout the paper we use \mathfrak{A} generically to represent one such family of ranges.

Also recall from Section 2 that if (P, \mathfrak{A}) has bounded VC-dimension ν , then we can create an ε -sample, with probability $1 - \delta$, by sampling $O(\frac{\nu}{\varepsilon^2} \log \frac{\nu}{\varepsilon \delta})$ points at random, or deterministically of size $O(\frac{\nu}{\varepsilon^2} \log \frac{\nu}{\varepsilon})$ in time $O(\nu^{2\nu} n^{(\frac{1}{\varepsilon^2} \log \frac{\nu}{\varepsilon})^\nu})$. There exist ε -samples of slightly smaller sizes [18], but efficient constructions are not known. If (P, \mathfrak{A}) has VC-dimension ν , this also implies that (P, \mathfrak{A}) contains at most $|P|^\nu$ sets.

Similarly, the *shatter function* $\pi_{(P, \mathfrak{A})}(m)$ of a range space (P, \mathfrak{A}) is the maximum number of sets $S \in (P, \mathfrak{A})$ where $|S| = m$. The *shatter dimension* σ of a range space (P, \mathfrak{A}) is the minimum value such that $\pi_{(P, \mathfrak{A})}(m) = O(m^\sigma)$. It can be shown [12] that $\sigma \leq \nu$ and $\nu = O(\sigma \log \sigma)$.

For a range space (P, \mathfrak{A}) the *dual range space* is defined (\mathfrak{A}, P^*) where P^* is all subsets $\mathfrak{A}_p \subseteq \mathfrak{A}$ defined for an element $p \in P$ such that $\mathfrak{A}_p = \{A \in \mathfrak{A} \mid p \in A\}$. If (P, \mathfrak{A}) has VC-dimension ν , then (\mathfrak{A}, P^*) has VC-dimension

$\leq 2^{\nu+1}$. Thus, if the VC-dimension of (\mathbb{A}, P^*) is constant, then the VC-dimension of (P, \mathbb{A}) is also constant [17]. Hence, the standard ε -sample theorems apply to dual range spaces as well.

Let $g : \mathbb{R} \rightarrow \mathbb{R}^+$ be a function where $\int_{x=-\infty}^{\infty} g(x) dx = 1$. We can create an ε -sample Q_g of (g, \mathbb{I}_+) , where \mathbb{I}_+ describes the set of all one-sided intervals of the form $(-\infty, t)$, so that

$$\max_t \left| \int_{x=-\infty}^t g(x) dx - \frac{1}{|Q_g|} \sum_{q \in Q_g} 1(q < t) \right| \leq \varepsilon.$$

We can construct Q_g of size $O(\frac{1}{\varepsilon})$ by choosing a set of points in Q_g so that the integral between two consecutive points is always ε . But we do not need to be so precise. Consider the set of $\frac{2}{\varepsilon}$ points $\{q'_1, q'_2, \dots, q'_2\}$ such that $\int_{x=-\infty}^{q'_i} = i\varepsilon/2$. Any set of $\frac{2}{\varepsilon}$ points $Q_g = \{q_1, q_2, \dots, q_2\}$ such that $q'_i \leq q_i \leq q'_{i+1}$ is an ε -sample.

A.1 ε -Samples of Distributions.

We say a subset $W \subset \mathbb{R}^d$ is *polygonal approximable* if there exists a polygonal shape S with m facets such that $\phi(W \setminus S) + \phi(S \setminus W) \leq \varepsilon\phi(W)$ for any $\varepsilon > 0$. Usually, m is dependent on ε . In turn, such a polygonal shape S describes a continuous point set where (S, \mathbb{A}) can be given an ε -sample Q using $O(\frac{1}{\varepsilon^2} \log \frac{1}{\varepsilon})$ points if (S, \mathbb{A}) has bounded VC-dimension [17] or using $O(\frac{1}{\varepsilon} \log^{2k} \frac{1}{\varepsilon})$ points if \mathbb{A} is defined by a constant k number of directions [21]. For instance, where $\mathbb{A} = \mathbb{B}$ is the set of all balls then the first case applies, and when $\mathbb{A} = \mathbb{R}_2$ is the set of all axis-aligned rectangles then either case applies.

A shape $W \subset \mathbb{R}^{d+1}$ may describe a distribution $\mu : \mathbb{R}^d \rightarrow [0, 1]$. We note that many common distributions like multivariate Gaussian distributions are polygonally approximable. For instance for a range space (μ, \mathbb{B}) , then the range space of the associated shape W_μ is $(W_\mu, \mathbb{B} \times \mathbb{R})$ where $\mathbb{B} \times \mathbb{R}$ describes balls in \mathbb{R}^d for the first d coordinates and any points in the $(d+1)$ th coordinate.

The general scheme to create an ε -sample for (S, \mathbb{A}) , where $S \in \mathbb{R}^d$ is a polygonal shape, is to use a lattice Λ of points. A *lattice* Λ in \mathbb{R}^d is an infinite set of points defined such that for d vectors $\{v_1, \dots, v_d\}$ that form a basis, for any point $p \in \Lambda$, $p + v_i$ and $p - v_i$ are also in Λ for any $i \in [1, d]$. We first create a discrete $\frac{\varepsilon}{2}$ -sample $M \subset \Lambda$ of (S, \mathbb{A}) and then create an $\frac{\varepsilon}{2}$ -sample Q of (M, \mathbb{A}) using standard techniques [5, 21]. Then Q is an ε -sample of (S, \mathbb{A}) . For a shape S with m $(d-1)$ -faces on its boundary, any subset $A' \subset \mathbb{R}^d$ that is described by a subset from (S, \mathbb{A}) is an intersection $A' = A \cap S$ for some $A \in \mathbb{A}$. Since S has m $(d-1)$ -dimensional faces, we can bound the VC-dimension of (S, \mathbb{A}) as $\nu = O((m + \nu_{\mathbb{A}}) \log(m + \nu_{\mathbb{A}}))$ where $\nu_{\mathbb{A}}$ is the VC-dimension of $(\mathbb{R}^d, \mathbb{A})$. Finally the set $M = S \cap \Lambda$ is determined by choosing an arbitrary initial origin point in Λ and then uniformly scaling all vectors $\{v_1, \dots, v_d\}$ until $|M| = \Theta(\frac{\nu}{\varepsilon^2} \log \frac{\nu}{\varepsilon})$ [17]. This construction follows a less general but smaller construction in Phillips [21].

It follows that we can create such an ε -sample of size $|M|$ in time $O(|M|m \log |M|)$ by starting with a scaling of the lattice so a constant number of points are in S and then doubling the scale until we get to within a factor of d of $|M|$. If there are n points inside S , it takes $O(nm)$ time to count them.

Lemma A.1. *For a polygonal shape $S \subset \mathbb{R}^d$ with m facets, we can construct an ε -sample for (S, \mathbb{A}) of size $O(\frac{\nu}{\varepsilon^2} \log \frac{\nu}{\varepsilon})$ in time $O(m \frac{\nu}{\varepsilon^2} \log^2 \frac{\nu}{\varepsilon})$, where (S, \mathbb{A}) has VC-dimension $\nu_{\mathbb{A}}$ and $\nu = O((\nu_{\mathbb{A}} + m) \log(\nu_{\mathbb{A}} + m))$.*

An important part of the above construction is the arbitrary choice of the origin points of the lattice Λ . This allows us to arbitrarily shift the lattice defining M and thus the set Q . In Section 6 we need to construct n ε -samples $\{Q_1, \dots, Q_n\}$ for n range spaces $\{(S_1, \mathbb{A}), \dots, (S_n, \mathbb{A})\}$. In Algorithm 6.1 we examine sets of $\nu_{\mathbb{A}}$ points, each from separate ε -samples that define a minimal shape $A \in \mathbb{A}$. It is important that we do not have two such (possibly not disjoint) sets of $\nu_{\mathbb{A}}$ points that define the same minimal shape $A \in \mathbb{A}$. (Note, this does not include cases where say two points are antipodal on a disk and any other point in the disk added to a set of $\nu_{\mathbb{A}} = 3$ points forms such a set; it refers to cases where say four points lie (degenerately) on the boundary of a disc.) We can guarantee this by enforcing a property on all pairs of origin points p and q for (S_i, \mathbb{A}) and (S_j, \mathbb{A}) . For the purpose of construction, it

is easiest to consider only the l th coordinates p_l and q_l for any pair of origin points or lattice vectors (where the same lattice vectors are used for each lattice). We enforce a specific property on every such pair p_l and q_l , for all l and all distributions and lattice vectors.

First, consider the case where $\mathbb{A} = \mathbb{R}_d$ describes axis-aligned bounding boxes. It is easy to see that if for all pairs p_l and q_l that $(p_l - q_l)$ is irrational, then we cannot have $> 2d$ points on the boundary of an axis-aligned bounding box, hence the desired property is satisfied.

Now consider the more complicated case where $\mathbb{A} = \mathbb{B}$ describes smallest enclosing balls. There is a polynomial of degree 2 that describes the boundary of the ball, so we can enforce that for all pairs p_l and q_l that $(p_l - q_l)$ is of the form $c_1(r_{p_l})^{1/3} + c_2(r_{q_l})^{1/3}$ where c_1 and c_2 are rational coefficients and r_{p_l} and r_{q_l} are distinct integers that are not multiple of cubes. Now if $\nu = d + 1$ such points satisfy (and in fact define) the equation of the boundary of a ball, then no $(d + 2)$ th point which has this property with respect to the first $d + 1$ can also satisfy this equation.

More generally, if \mathbb{A} can be described with a polynomial of degree p with ν variables, then enforce that every pair of coordinates are the sum of $(p + 1)$ -roots. This ensures that no $\nu + 1$ points can satisfy the equation, and the undesired situation cannot occur.

B A Center Point for μ_P

We can create a point $\bar{q} \in \mathbb{R}^d$ that is in the convex hull of a sampled point set Q from μ_P with high probability. This implies that for any summarizing shape that contains the convex hull, \bar{q} is also contained in that summarizing shape. Let \mathbb{H} be the family of subsets defined by halfspaces. We use the following algorithm:

1. Create 2-approximate center points \bar{p}_i for each μ_{p_i} (i.e. using a $(1/4)$ -sample of (μ_{p_i}, \mathbb{H})). Let the set be \bar{P} .
2. Create 2-approximate center point \bar{q} of \bar{P} .

All steps can be done in $O(n)$ time because we can create $(1/4)$ -samples of all range spaces (μ_{p_i}, \mathbb{H}) and of (\bar{P}, \mathbb{H}) in $O(n)$ time. Constructing approximate center points can be done in $O(1)$ time on a constant sized set, such as $(1/4)$ -sample [7].

Lemma B.1. *Given a distribution of a point set μ_P (such that each point distribution is polygonally approximable) of n points in \mathbb{R}^d , there is an $O(n)$ time algorithm to create a point \bar{q} that will be in the convex hull of a point set drawn from μ_P with probability $\geq 1 - ((1 - 1/(2d + 2))^{1/(2d+2)})^n$.*

Proof. For each $\bar{p}_i \in \bar{P}$, any halfspace that has \bar{p}_i on its boundary and does not contain \bar{q} has probability $\geq 1/(2d+2)$ of containing a random point from μ_{p_i} . Thus for any direction $u \in \mathbb{S}^{d-1}$ there are at least $n/(2d + 2)$ points \bar{p}_i from \bar{P} for which $\langle \bar{q}, u \rangle \leq \langle \bar{p}_i, u \rangle$. And for each of those points \bar{p}_i , the probability that the point q_i sampled from μ_{p_i} is such that $\langle \bar{p}_i, u \rangle \leq \langle q_i, u \rangle$ (and thus $\langle \bar{q}, u \rangle \leq \langle q_i, u \rangle$) is $\leq 1/(2d + 2)$. Hence, the probability that there is a separating halfspace between \bar{q} and the convex hull of Q (where the halfspace is orthogonal to some direction u) is $\leq (1 - 1/(2d + 2))^{n/(2d+2)} = ((1 - 1/(2d + 2))^{1/(2d+2)})^n$. \square

Theorem B.1. *For a set of $m < n$ point sets drawn i.i.d. from μ_P , it follows that \bar{q} is in each of the m convex hulls for each point sets with high probability (specifically with probability $\geq 1 - m ((1 - 1/(2d + 2))^{1/(2d+2)})^n$).*

Proof. Let $\beta = (1 - 1/(2d + 2))^{1/(2d+2)}$. For any one point set the probability that \bar{q} is contained in the convex hull is $> 1 - \beta^n$. By the union bound, the probability that it is contained in all m convex hulls is $> (1 - \beta^n)^m = 1 - m\beta^n + \binom{m}{2}\beta^{2n} - \binom{m}{3}\beta^{3n} + \dots$. Since $n > m$, the sum of all terms after the first two in the expansion increase the probability. \square

Thus because the summarizing shapes are convex, then for any point q , the line segment $\overline{q\bar{q}}$ is completely contained in a convex summarizing shape if and only if q is. Thus for every boundary of a summarizing shape $\overline{q\bar{q}}$ crosses, q is outside that summarizing shape. This implies the following corollary.

Corollary B.1. *If the summarizing shape is convex, then the γ -layer, for $\gamma < 1 - 1/m$, exists, is connected, and is star-shaped with high probability, specifically with probability $\geq 1 - m ((1 - 1/(2d + 2))^{1/(2d+2)})^n$.*

C Shapes of $\mathbb{A}_{f,n}$ for Various Summarizing Shapes

Let $\mathbb{A}_{f,n}$ be the intersection of $O(n)$ shapes from $\hat{\mathbb{A}}_f$ where $(\mu_p, \hat{\mathbb{A}}_f)$ has VC-dimension $\hat{\nu}_f$. Let a *wedge* of $\mathbb{A}_{f,n}$ be a shape from $\mathbb{W}_{f,n}$ described by the intersection of d hyperplanes and one shape from $\hat{\mathbb{A}}_f$.

Lemma C.1. *The VC-dimension of $(\mathbb{R}^d, \mathbb{W}_{f,n})$ is $d(d+1) + \hat{\nu}_f$.*

Proof. It is known that a class of shapes \mathbb{W}_k that is formed as the intersection of k subset from $(\mathbb{R}^d, \mathbb{A}_j)$ for $j \in [1 : k]$ which have VC-dimension ν_j , then $(\mathbb{R}^d, \mathbb{W}_k)$ has VC-dimension $\sum_{j=1}^k \nu_j$ [12]. Since wedges are formed by the intersection of d halfspaces (VC-dimension $d+1$) and one shape from $\hat{\mathbb{A}}_f$, it follows that $(\mathbb{R}^d, \mathbb{W}_{f,n})$ has VC-dimension $d(d+1) + \hat{\nu}_f$. \square

If the $(d-2)$ -dimensional faces of the boundary of $\mathbb{A}_{f,n}$ are subsets of $(d-2)$ -dimensional flats (i.e. points in \mathbb{R}^2 and line segments in \mathbb{R}^3), then any shape from $\mathbb{A}_{f,n}$ can be formed as the disjoint union of $O(n)$ wedges from $\mathbb{W}_{f,n}$. Functions which produce such families $\mathbb{A}_{f,n}$ include seb_2 and chp in \mathbb{R}^2 and diam and cha in \mathbb{R}^d .

Remark 1. In cases, such as seb_2 and chp for $d > 2$, where we cannot form wedges, we can create similar shapes for $\mathbb{W}_{f,n}$, as generalized cones whose boundary passes through the boundary of each $(d-1)$ -dimensional facet of the corresponding shape from $\mathbb{A}_{f,n}$. For these shapes the VC-dimension of $(\mathbb{R}^d, \mathbb{W}_{f,n})$ can be bounded as $O(\nu_f d \log d)$. Each face of a generalized cone is described by two shapes from \mathbb{A} , which have VC-dimension ν_f , and a point. Thus the face of the generalized cone has shatter dimension $O(\nu_f)$. If a $(d-1)$ -dimensional facet of the boundary of a shape from $\mathbb{A}_{f,n}$ has more than d faces of dimension $(d-2)$, then we can triangulate the facet so it has $O(d)$ such faces. Thus the range space for the generalized cone has shatter dimension $O(\nu_f d)$ and VC-dimension $O(\nu_f d \log d)$.

The VC-dimension for $(\mathbb{R}^d, \mathbb{W}_{f,n})$, ψ_f , is shown for several functions in Table 2.

Lemma C.2. *If the disjoint union of m shapes from $\mathbb{W}_{f,n}$ can form any shape from $\mathbb{A}_{f,n}$, then an $\frac{\varepsilon}{m}$ -sample of $(\mu_p, \mathbb{W}_{f,n})$ is an ε -sample of $(\mu_p, \mathbb{A}_{f,n})$.*

Proof. For any shape $A \in \mathbb{A}_{f,n}$ we can create a set of m shapes $\{W_1, \dots, W_m\} \subset \mathbb{W}_{f,n}$ whose disjoint union is A . Since each range of $\mathbb{W}_{f,n}$ may have error $\frac{\varepsilon}{m}$, their union has error at most ε . \square

Hence for $(\mu_p, \mathbb{W}_{f,n})$ $\frac{\varepsilon}{n}$ -samples can be created of size $O(n^2 \frac{1}{\varepsilon^2} \log \frac{n}{\varepsilon})$ in time $O(n(\frac{n}{\varepsilon})^{2\psi_f} \log^{\psi_f} \frac{n}{\varepsilon})$.

Table 1: Runtimes for ε -Quantizations of Various Summarizing Shape Families.

abbrv.	summarizing shape	randomized*	determ. \mathbb{R}^2	determ. \mathbb{R}^d
dwid	width along a fixed direction	$O(n \frac{1}{\varepsilon^2} \log \frac{1}{\varepsilon})$	$O(n^4/\varepsilon)$	$O(n^4/\varepsilon)$
aabbp	axis-aligned bounding box measured by perimeter	$O(n \frac{1}{\varepsilon^2} \log \frac{1}{\varepsilon})$	$\tilde{O}(n^8/\varepsilon^4)$	$\tilde{O}(n^{6d}/\varepsilon^{4d})$
aabba	axis-aligned bounding box measured by area	$O(n \frac{1}{\varepsilon^2} \log \frac{1}{\varepsilon})$	$\tilde{O}(n^{12}/\varepsilon^8)$	$\tilde{O}(n^{6d}/\varepsilon^{4d})$
seb $_\infty$	smallest enclosing ball, L_∞ metric	$O(n \frac{1}{\varepsilon^2} \log \frac{1}{\varepsilon})$	$\tilde{O}(n^6/\varepsilon^3)$	$\tilde{O}(n^{2d+2}/\varepsilon^{d+1})$
seb $_1$	smallest enclosing ball, L_1 metric	$O(n \frac{1}{\varepsilon^2} \log \frac{1}{\varepsilon})$	$\tilde{O}(n^6/\varepsilon^3)$	$\tilde{O}(n^{2d+2}/\varepsilon^{d+1})$
seb $_2$	smallest enclosing ball, L_2 metric	$O(n \frac{1}{\varepsilon^2} \log \frac{1}{\varepsilon})$	$\tilde{O}(n^{16.5}/\varepsilon^7)$	$\tilde{O}((n^{d+3}/\varepsilon^2)^{d+2-1/d})$
diam	diameter	$O(n^2 \frac{1}{\varepsilon^2} \log \frac{1}{\varepsilon})$	$\tilde{O}((n^5/\varepsilon^2)^{n+1})$	$\tilde{O}((n^5/\varepsilon^2)^{n+1})$
cha	convex hull measured by area	$O(n \log n \frac{1}{\varepsilon^2} \log \frac{1}{\varepsilon})$	$\tilde{O}((n^5/\varepsilon^2)^{n+1})$	$\tilde{O}((n^{d+3}/\varepsilon^2)^{n+1})$
chp	convex hull measured by perimeter	$O(n \log n \frac{1}{\varepsilon^2} \log \frac{1}{\varepsilon})$	$\tilde{O}((n^5/\varepsilon^2)^{n+1})$	$\tilde{O}((n^{d+3}/\varepsilon^2)^{n+1})$

* all randomized results are correct with constant probability.

$\tilde{O}(f(n, \varepsilon))$ ignores poly-logarithmic factors $(\log \frac{n}{\varepsilon})^{O(\text{poly}(d))}$, for any $\tau > 0$.

Table 2: VC-dimension for Various Shape Families.

abbrv.	$(\mathbb{R}^2, \mathbb{A}_{f,n})$	$(\mathbb{R}^d, \mathbb{A}_{f,n})$	$(\mathbb{R}^2, \mathbb{W}_{f,n})$	$(\mathbb{R}^d, \mathbb{W}_{f,n})$
aabbp	$O(1)$ ($\sigma = 4$)	$O(d \log d)$ ($\sigma = 2d$)		
aabba	8	$O(d \log d)$ ($\sigma = 2d$)		
seb $_{\infty}$	4	$2d$		
seb $_1$	4	$2d$		
seb $_2$	∞	∞	9	$O(d^2 \log d)$
diam	∞	∞	9	$d^2 + 2d + 1$
cha	∞	∞	7	$d^2 + 2d + 1$
chp	∞	∞	$O(1)$	$O(d^2 \log d)$

C.1 Examples

In the examples below, 7 points are given, on which we study a certain measure (e.g., diameter or convex hull area). The grey region denotes the possible placements of a new point, such that the measure will not exceed a given value. These regions illustrate $\mathbb{A}_{f,n}$ for various summarizing shapes.

Axis-aligned bounding box. Figure 5 shows examples of $\mathbb{A}_{f,n}$ for axis-aligned bounding boxes, measuring either by perimeter (aabbp) or by area (aabba) in \mathbb{R}^2 . For both $(\mathbb{R}^2, \mathbb{A}_{f,n})$ has a shatter dimension of 4 because the shape is determined by the x -coordinates of 2 points and the y -coordinates of 2 points. This generalizes to a shatter dimension of $2d$ for $(\mathbb{R}^d, \mathbb{A}_{f,n})$, where area generalizes to d -dimensional volume, and perimeter generalizes to the $(d-1)$ -volume of the boundary. We can also show the VC-dimension of $(\mathbb{R}^2, \mathbb{A}_{f,n})$ is 8 for aabbp because its shape is defined by the intersection of halfspaces with 4 predefined normal directions at 0° , 45° , 90° , and 135° . This can be generalized to higher dimensions.

Hence, for both shapes we can create $n \frac{\varepsilon}{n}$ -samples of $(\mu_{p_i}, \mathbb{A}_{f,n})$ of size $\alpha_f(n, \varepsilon) = O(\frac{n^2}{\varepsilon^2} \log \frac{n}{\varepsilon})$ in time $O(\frac{n^3}{\varepsilon^2} \log^2 \frac{n}{\varepsilon})$. For aabbp in \mathbb{R}^2 , an $\frac{\varepsilon}{n}$ -sample of each $(\mu_{p_i}, \mathbb{A}_{f,n})$ of can be reduced further to size $O(\frac{n}{\varepsilon} \log^{16} \frac{n}{\varepsilon})$ in total time $O(\frac{n^5}{\varepsilon^4} \log^{40} \frac{n}{\varepsilon})$. Then we can construct the ε -quantization in $(n^{6d}/\varepsilon^{4d})(\log \frac{n}{\varepsilon})^{O(d)}$ time, using orthogonal range searching. For aabba in \mathbb{R}^2 , the runtime improves to $O(n^8/\varepsilon^4 \log^{65} \frac{n}{\varepsilon})$.

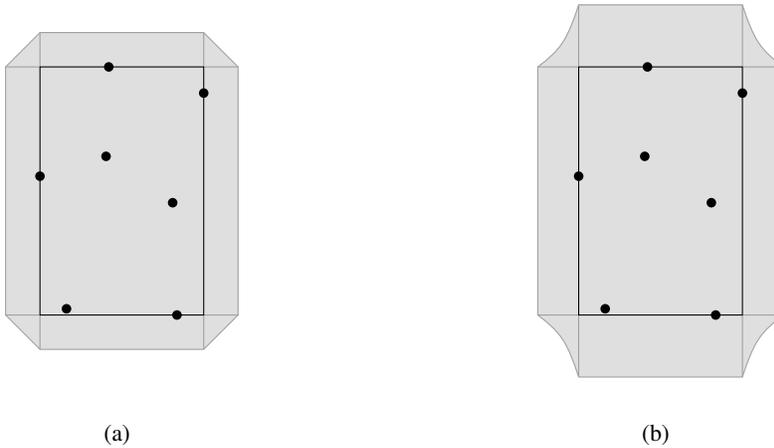


Figure 5: (a) Axis-aligned bounding box, measured by perimeter. (b) Axis-aligned bounding box, measured by area. The curves are hyperbola parts.

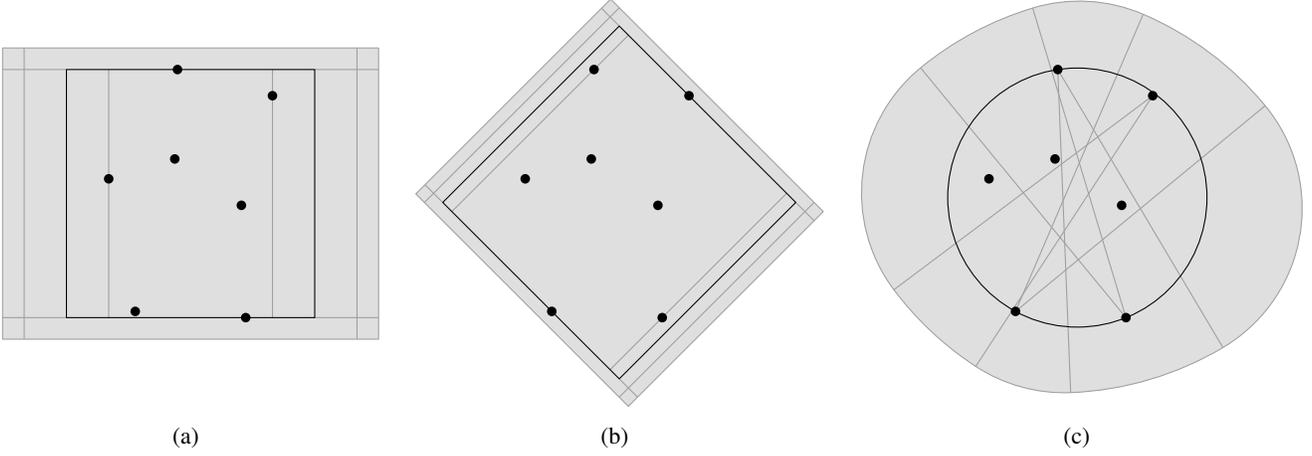


Figure 6: (a) Smallest enclosing ball, L_∞ metric. (b) Smallest enclosing ball, L_1 metric. (c) Smallest enclosing ball, L_2 metric. The curves are circular arcs of two different radii.

Smallest enclosing ball. Figure 6 shows examples of $\mathbb{A}_{f,n}$ for smallest enclosing ball, for metrics L_∞ (seb_∞), L_1 (seb_1), and L_2 (seb_2) in \mathbb{R}^2 . For seb_∞ and seb_1 , $(\mathbb{R}^d, \mathbb{A}_{f,n})$ has VC-dimension $2d$ because the shapes are defined by the intersection of halfspaces from d predefined normal directions. For seb_1 and seb_∞ , we can create $n \frac{\varepsilon}{n}$ -samples of each $(\mu_{p_i}, \mathbb{A}_{f,n})$ of size $\alpha_f(n, \varepsilon) = O(\frac{n^2}{\varepsilon^2} \log \frac{n}{\varepsilon})$ in total time $O(\frac{n^3}{\varepsilon^2} \log^2 \frac{n}{\varepsilon})$. The size for each can be reduced to $O(\frac{n}{\varepsilon} \log^{2d} \frac{n}{\varepsilon})$ in $O(\frac{n^5}{\varepsilon^4} \log^{8d} \frac{n}{\varepsilon})$ total time. Using an orthogonal range searching data structure we can calculate the ε -quantization in $O(n^{2d+2}/\varepsilon^{d+1} \log^{7d-1} \frac{n}{\varepsilon})$ time.

For seb_2 , $(\mathbb{R}^d, \mathbb{A}_{f,n})$ has infinite VC-dimension, but $(\mathbb{R}^2, \mathbb{W}_{f,n})$ has VC-dimension ≤ 9 because it is the intersection of 2 halfspaces and one disc. Any shape from $\mathbb{A}_{f,n}$ can be formed from the disjoint union of $2n$ wedges. Choosing a point in the convex hull of the $n - 1$ points describing $\mathbb{A}_{f,n}$ as the vertex of the wedges will ensure that each wedge is completely inside the ball that defines part of its boundary. Thus, in \mathbb{R}^2 the $n \frac{\varepsilon}{n}$ -samples of each $(\mu_{p_i}, \mathbb{A}_{f,n})$ are of size $\alpha_f(n, \varepsilon) = O(n^4/\varepsilon^2 \log \frac{n}{\varepsilon})$ and can all be calculated in $O(n^5/\varepsilon^2 \log^2 \frac{n}{\varepsilon})$ time. And then the ε -quantization can be calculated in $O(n^{16.5}/\varepsilon^7 \log^{3.5} \frac{n}{\varepsilon})$ time, using range searching data structures.

For seb_2 , in \mathbb{R}^d for $d > 3$, we can form shapes from $\mathbb{A}_{f,n}$ with disjoint unions of generalized cone from a family $\mathbb{W}_{f,n}$, where $(\mathbb{R}^d, \mathbb{W}_{f,n})$ has shatter dimension $O(d^2)$. We need $O(n^{\lfloor d/2 \rfloor})$ such shapes from $\mathbb{W}_{f,n}$ to form one shape $A \in \mathbb{A}_{f,n}$, because A has boundary described by $O(n^{\lfloor d/2 \rfloor})$ sphere piece with one of d different radii. The VC-dimension of each $(\mu_{p_i}, \mathbb{W}_{f,n})$ is $O(d^2 \log d)$ in \mathbb{R}^d , and we can create $n \frac{\varepsilon}{n}$ -quantization of each $(\mu_{p_i}, \mathbb{A}_{f,n})$ of size $\alpha_f(n, \varepsilon) = O(n^{2+2\lfloor d/2 \rfloor}/\varepsilon^2 \log \frac{n}{\varepsilon})$ in $O(n^{3+2\lfloor d/2 \rfloor}/\varepsilon^2 \log^2 \frac{n}{\varepsilon})$ total time. Then the ε -quantization can be computed in $O((n^{d+3}/\varepsilon^2)^{d+2-1/d})$ time (ignoring boundary cases with floor operations) using a range searching data structure.

Diameter. Figure 7 shows an example of $\mathbb{A}_{f,n}$ for the diameter of a point set in \mathbb{R}^2 . Here $(\mathbb{R}^d, \mathbb{A}_{f,n})$ has infinite VC-dimension. It is formed by the intersection of balls of the same radius centered at the points. Thus a shape from $\mathbb{A}_{f,n}$ is determined by at most n balls, and since they are each the same radius, we can construct a shape from $\mathbb{A}_{f,n}$ from the disjoint union of n wedges, as with seb_2 . And since each wedge is the intersection of d halfspaces and 1 disc, $(\mathbb{R}^d, \mathbb{W}_{f,n})$ has VC-dimension $(d + 1)^2$. Thus we can construct $n \frac{\varepsilon}{n}$ -samples for each $(\mu_{p_i}, \mathbb{A}_{f,n})$ of size $\alpha_f(n, \varepsilon) = O(n^4/\varepsilon^2 \log \frac{n}{\varepsilon})$ in total time $O(n^5/\varepsilon^2 \log^2 \frac{n}{\varepsilon})$. However, given a set of n points, the shape which defines the set of points where the diameter will not increase has complexity $O(n)$. This is the family $\mathbb{A}_{f,n+1}$, the union of n balls. This implies the size of the basis σ_f used in Algorithm 6.1 is n . Hence, it takes time $O((n^5/\varepsilon^2 \log \frac{n}{\varepsilon})^{n+1})$ to construct the ε -quantization.

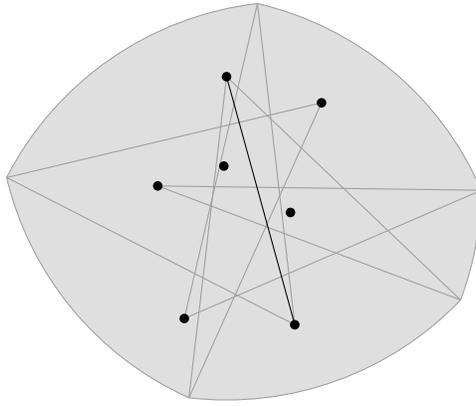


Figure 7: Diameter. The curves are circular arcs all of the same radius.

Convex hull. Figure 8 shows examples of $\mathbb{A}_{f,n}$ for the convex hull, measured either by area (cha) or perimeter (chp) in \mathbb{R}^2 . For both, $(\mathbb{R}^2, \mathbb{A}_{f,n})$ has infinite VC-dimension. For cha $(\mathbb{R}^2, \mathbb{W}_{f,n})$ has VC-dimension 7, because wedges are triangles. In higher dimensions cha can continue to use wedges, but needs $O(n^{\lfloor d/2 \rfloor})$ of them. For chp, the wedges boundary is described by d hyperplanes and an ellipse boundary part. We cannot guarantee that the intersection of all of these parts describes the wedge because the ellipse may be too small and may cut off part of the intersection of halfspaces. But in \mathbb{R}^2 the wedge clearly does have shatter dimension $4 + 5$, so the VC-dimension of $(\mathbb{R}^2, \mathbb{W}_{f,n})$ is $O(1)$. In higher dimensions we can use generalized cone shapes with VC-dimension $O(d^2 \log d)$ and we may need $O(n^{\lfloor d/2 \rfloor})$ of them.

For both cha and chp we can calculate $n \frac{\varepsilon}{n}$ -samples for each $(\mu_{p_i}, \mathbb{A}_{f,n})$ of size $\alpha_f(n, \varepsilon) = O(n^{2+2\lfloor d/2 \rfloor} / \varepsilon^2 \log \frac{n}{\varepsilon})$ in $O(n^{3+2\lfloor d/2 \rfloor} / \varepsilon^2 \log^2 \frac{n}{\varepsilon})$ total time. However, given a set of n points, the shape which defines the set of points where the convex hull will not increase has complexity $O(n)$. This is the family $\mathbb{A}_{f,n+1}$. Like diam, this implies the size of the basis σ_f used in Algorithm 6.1 is n . Hence, it takes time $O((n^{d+3} / \varepsilon^2 \log \frac{n}{\varepsilon})^{n+1})$ to construct the ε -quantizations.

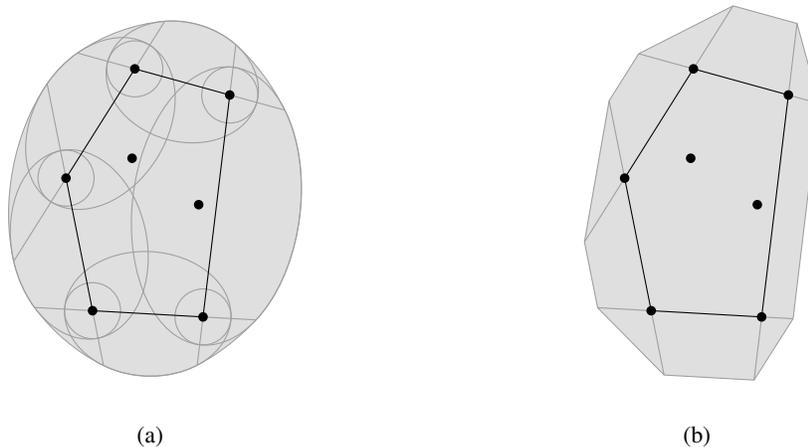


Figure 8: (a) Convex hull, measured by perimeter. The curves are ellipse parts. (b) Convex hull, measured by area.

Table 3: ε -Samples for Summarizing Shape Family $\mathbb{A}_{f,n}$.

abbrv.	$\alpha_f(n, \varepsilon)$	$\eta = n\alpha_f(n, \varepsilon)$	ν_f	η^{ν_f}	$\text{RS}_f(n, \varepsilon)$	runtime
aabbp	$\tilde{O}(\frac{n^2}{\varepsilon^2})$	$\tilde{O}(n^3/\varepsilon^2)$	$2d$	$\tilde{O}(n^{6d}/\varepsilon^{4d})$	$\tilde{O}(1)$	$\tilde{O}(n^{6d}/\varepsilon^{4d})$
aabba	$\tilde{O}(\frac{n^2}{\varepsilon^2})$	$\tilde{O}(n^3/\varepsilon^2)$	$2d$	$\tilde{O}(n^{6d}/\varepsilon^{4d})$	$\tilde{O}(1)$	$\tilde{O}(n^{6d}/\varepsilon^{4d})$
seb $_{\infty}$	$\tilde{O}(\frac{n}{\varepsilon})$	$\tilde{O}(n^2/\varepsilon)$	$d+1$	$\tilde{O}(n^{2d+2}/\varepsilon^{d+1})$	$\tilde{O}(1)$	$\tilde{O}(n^{2d+2}/\varepsilon^{d+1})$
seb $_1$	$\tilde{O}(\frac{n}{\varepsilon})$	$\tilde{O}(n^2/\varepsilon)$	$d+1$	$\tilde{O}(n^{2d+2}/\varepsilon^{d+1})$	$\tilde{O}(1)$	$\tilde{O}(n^{2d+2}/\varepsilon^{d+1})$
seb $_2$	$\tilde{O}(\frac{n^{d+2}}{\varepsilon^2})$	$\tilde{O}(n^{d+3}/\varepsilon^2)$	$d+1$	$\tilde{O}((n^{d+3}/\varepsilon^2)^{d+1})$	$O((n^{d+3}/\varepsilon^2)^{1-1/d})$	$\tilde{O}((n^{d+3}/\varepsilon^2)^{d+2-1/d})$
diam	$\tilde{O}(\frac{n^4}{\varepsilon^2})$	$\tilde{O}(n^5/\varepsilon^2)$	n	$\tilde{O}(n^5/\varepsilon^2)^n$	$O(n^5/\varepsilon^2)$	$\tilde{O}((n^5/\varepsilon^2)^{n+1})$
cha	$\tilde{O}(\frac{n^{d+2}}{\varepsilon^2})$	$\tilde{O}(n^{d+3}/\varepsilon^2)$	n	$\tilde{O}((n^{d+3}/\varepsilon^2)^n)$	$O(n^{d+3}/\varepsilon^2)$	$\tilde{O}((n^{d+3}/\varepsilon^2)^{n+1})$
chp	$\tilde{O}(\frac{n^{d+2}}{\varepsilon^2})$	$\tilde{O}(n^{d+3}/\varepsilon^2)$	n	$\tilde{O}((n^{d+3}/\varepsilon^2)^n)$	$O(n^{d+3}/\varepsilon^2)$	$\tilde{O}((n^{d+3}/\varepsilon^2)^{n+1})$

$\tilde{O}(f(n, \varepsilon))$ ignores poly-logarithmic factors $(\log \frac{n}{\varepsilon})^{O(\text{poly}(d))}$.