



ELSEVIER

Theoretical Computer Science 287 (2002) 501–514

Theoretical
Computer Science

www.elsevier.com/locate/tcs

On computing the diameter of a point set in high dimensional Euclidean space[☆]

Daniele V. Finocchiaro^a, Marco Pellegrini^{b,*}

^a*Scuola Normale Superiore in Pisa, Piazza dei Cavalieri 7, I-56126 Pisa, Italy*

^b*Institute for Computational Mathematics of CNR via S. Maria 46, I-56126 Pisa, Italy*

Received March 2001; accepted May 2001

Abstract

We consider the problem of computing the diameter of a set of n points in d -dimensional Euclidean space under Euclidean distance function. We describe an algorithm that in time $O(dn \log n + n^2)$ finds with high probability an arbitrarily close approximation of the diameter. For large values of d the complexity bound of our algorithm is a substantial improvement over the complexity bounds of previously known exact algorithms. Computing and approximating the diameter are fundamental primitives in high dimensional computational geometry and find practical application, for example, in clustering operations for image databases. © 2002 Elsevier Science B.V. All rights reserved.

1. Introduction

We consider the following problem: given a set of n points in d -dimensional Euclidean space, compute the maximum pair-wise Euclidean distance between two points in the set. This problem is known as the *diameter problem* or the *furthest pair problem*. There are several efficient algorithms for the case when $d = 2$ [31] and $d = 3$ [3, 32, 8, 7, 27] which, however, do not extend to higher dimensional spaces. The diameter problem is one of the basic problems in high dimensional computational geometry [17–19].¹ In this paper, we consider a setting in which the number of points n and the dimension of the space d are equally important in the complexity analysis.

[☆] A preliminary version of this paper has appeared in the Proceedings of the 7th Annual European Symposium on Algorithms (ESA '99) pp. 366–377.

* Corresponding author. Area della Ricerca, IMC-CNR, Via Moruzzi 1 I-56124 Pisa, Italy.
E-mail addresses: fino@cibs.sns.it (D.V. Finocchiaro), pellegrini@imc.pi.cnr.it (M. Pellegrini).

¹ Gritzmann and Klee have proposed the term *Computational Convexity* to denote the study of combinatorial and algorithmic aspects of polytopes in high dimensional spaces.

The exact solution to the diameter problem in arbitrary dimension can be found using the trivial algorithm that generates all possible $\Theta(n^2)$ inter-point distances and determines the maximum value. This algorithm runs in time $O(dn^2)$. As noted in [31] substantial improvements of the asymptotic complexity in terms of n and d must overcome the fact that for $d \geq 4$ the number of diametral pairs of points can be $\Omega(n^2)$ [14],² while computing a single inter-point distance takes time $\Omega(d)$. Yao in [38] gives an algorithm to compute the diameter in time $O(n^{2-a(d)} \log^{1-a(d)} n)$ where $a(d) = 2^{-(d+1)}$ for $d \geq 3$ and fixed. The technique in [38] can be extended to an $o(n^2)$ algorithm in non-fixed dimension for $d \leq \frac{1}{2} \log \log n$.

A result of Yao [37] (cited in [38]) shows that all inter-point distances can be computed in time $O(M(n, d) + nd + n^2)$, where $M(n, d)$ is the time to multiply an $n \times d$ matrix by a $d \times n$ matrix. Using the asymptotically fastest known square matrix multiplication algorithm [9] we have $M(n, d) \leq O(n^2 d^{s-2})$ when $n > d$, and $M(n, d) \leq O(dn^{s-1})$ when $n < d$, where $s \approx 2.376$. The furthest pair is then found trivially with $O(n^2)$ extra time.

In order to obtain faster algorithms for large values of d , we relax the requirement by considering algorithms that approximate the diameter up to multiplicative factors arbitrarily close to one. A powerful tool in this approximate setting is the Johnson–Lindenstrauss (JL) Lemma (see [24, 16]), showing that a projection of the data point set on a random linear subspace of dimension $O(\log n)$ preserves with high probability all inter-point distances within a multiplicative factor arbitrarily close to 1. Thus, with an overhead cost of $O(nd \log n)$ operations we can reduce the dimension of the ambient space to $d' \leq O(\log n)$, when $d \geq \Omega(\log n)$. Afterwards, using the trivial exact algorithm, we attain a bound $O(\min\{nd \log n + n^2 \log n, dn^2\})$. If, instead, Yao's second algorithm is used it is possible to achieve a bound $O(\min\{nd \log n + n^2 \log^{s-2} n, n^2 d^{s-2}\})$.

The main result in this paper (Theorem 3) is the following: the diameter of a point-set of size n in dimension d is approximated in time $O(dn \log n + n^2)$ within a factor $1 + \varepsilon$, for any real value $\varepsilon > 0$, with probability $1 - \delta$. The constants hidden in the big-Oh notation depend on user controlled parameters ε and δ , but not on d .

Our result improves asymptotically, by a factor $\min\{d, \log n\}$, over the trivial exact algorithm with the LJ projection in the range $\Omega(1) \leq d \leq n$. Our algorithm improves asymptotically, by a factor $(\min\{d, \log n\})^{s-2}$, over Yao's second algorithm with the LJ projection for d in the range: $\frac{1}{2} \log \log n \leq d \leq n$. For $d > n$ the overhead due to the application of the LJ lemma dominates any other cost in all the above mentioned algorithms.

Another approximation algorithm in literature, [12], attains a fixed approximation factor $c = \sqrt{5 - 2\sqrt{3}}$, which is not arbitrarily close to one.

In a recent paper Borodin et al. [5], amongst other important results, solve the approximate furthest pair problem in time roughly $O(nd \log n + n^{2-\theta(\varepsilon^2)})$ with high probability. Their result uses the LJ projection and improves asymptotically over our

² Instead, in dimensions 2 and 3, the number of diametral pairs of points is $O(n)$ [13, 33].

bound for $d < n / \log n$, however, for reasonable values of ε , say $\varepsilon = 0.1$, the asymptotic improvement is of the order of $n^{1/100}$, and such gain should be weighted against a more complex coding effort.

1.1. Applications

One of the most interesting applications of the computation of the diameter is in *clustering*: we are given a set of objects, and we want to group them in such a way that objects in the same group are “similar” [21]. One strategy is to map the objects to points in Euclidean space, so that similarity is mapped into closeness, and to seek for clusters of points. One of the most natural measure of quality for a cluster is its diameter. One application is data clustering for images database [20]. To map an image into a point in an high dimensional space, we associate a dimension to each term of a wavelet expansion of the image considered as a two-dimensional piecewise constant function [15, 23]. Thus similarity clustering of images is translated in the problem of determining clusters of such points. Other applications of clustering algorithms are in statistics, pattern recognition [30], biology, web search engines [6], distributed networks. Note that in such applications the number of dimensions is very high, thus the complexity due to the dimension must be taken into consideration in the design of efficient algorithms. On the other hand, clustering algorithms are usually based on heuristic arguments and approximations, so using an approximation of the diameter is as good as using its exact value.

A different application of the diameter computation is in a greedy heuristic for the maximum-weight matching problem on Euclidean graphs [2, 34]: the greedy algorithm iteratively matches a farthest pair, and deletes those two points from the set, until there are no more points. The weight of the matching produced is at least half of the weight of the maximum matching.

1.2. The algorithm

Our algorithm has been inspired by a recent technique for nearest neighbor search described by Kleinberg [25]. Although the formulation of the closest pair problem seems not too different from that of the diameter (searching for the smallest inter-point distance instead of the maximum), the mathematical properties of the two problems are quite different so that in general not any efficient algorithm for the closest pair problem yields an efficient one for the farthest pair problem.

Intuitively, the method of Kleinberg is based on the idea that if a vector $x \in \mathbb{R}^d$ is longer than a vector $y \in \mathbb{R}^d$ then this relation is preserved with probability greater than $\frac{1}{2}$ in a projection of x and y over a random line. Thus using several projections of the same set of vectors and a majority voting scheme we can retrieve, with probability close to one, the “actual” relative length relation between x and y . The theory of range spaces of bounded VC-dimension is invoked to determine how many random lines are needed to satisfy the constraints imposed by the error parameter ε and the confidence parameter δ .

1.3. Organization of the paper

In Section 2, we review some basic properties of projections of random vectors and of range spaces with bounded VC-dimension. In Section 3, we give the preprocessing and query algorithms for furthest point queries. Finally in Section 4, we apply the data structures of Section 3 to the problem of determining the diameter, thus establishing the main result.

2. Basic results

We denote with S^{d-1} the set of directions in \mathbb{R}^d ; it can be identified with the set of unit vectors, or with the set of points on the unit sphere. For any two vectors $x, y \in \mathbb{R}^d$, we define the set of directions over which the length of the projection of x is longer or equal than that of y :

$$\mathcal{L}_{x,y} = \{v \in S^{d-1} : |v \cdot x| \geq |v \cdot y|\}.$$

We denote with $z_{x,y}$ the relative measure of $\mathcal{L}_{x,y}$ with respect to the entire set of directions; it represents the probability that the projection of x on a random direction is longer than the projection of y :

$$z_{x,y} = \frac{|\mathcal{L}_{x,y}|}{|S^{d-1}|} = \Pr[|v \cdot x| \geq |v \cdot y|].$$

Note that for $x \neq y$ the probability that $|v \cdot x| = |v \cdot y|$ is null. Fig. 1 shows, by means of simple geometric intersections, the angles *orthogonal* to directions in $\mathcal{L}_{x,y}$ (marked as “good”).

The idea of this work is that if x is ‘significantly’ longer than y , then the set $\mathcal{L}_{x,y}$ is large, and it is very probable that a random direction belongs to it. The following lemma captures this idea (for a complete proof see [25]; here we detail some minor corrections):

Lemma 1. *If $(1 - \gamma)\|x\| \geq \|y\|$, with $0 \leq \gamma \leq 1$, then $z_{x,y} \geq \frac{1}{2} + \gamma/\pi$.*

Proof. Consider the plane \mathcal{A} containing the vectors x, y , and the projection $\pi_{\mathcal{A}}(v)$ of v on this plane. Observe that $v \cdot x = \pi_{\mathcal{A}}(v) \cdot x$ and $v \cdot y = \pi_{\mathcal{A}}(v) \cdot y$. So the relation

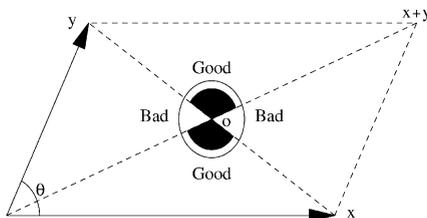


Fig. 1. Good and bad projection angles.

$v \in \mathcal{L}_{x,y}$ depends only on $\pi_{\mathcal{A}}(v)$, or equivalently on the angle ϕ between $\pi_{\mathcal{A}}(v)$ and the vector x (or y).

Note that ϕ is uniformly distributed in $[0, 2\pi)$. Let θ be the angle formed by x and y : we find that $v \in \mathcal{L}_{x,y}$ whenever $\cos^2(\theta - \phi) \leq \cos^2 \phi \cdot \|x\|^2 / \|y\|^2$, and the quantity $z_{x,y}$ is minimum when x and y are orthogonal. In this case

$$z_{x,y} = \frac{2}{\pi} \tan^{-1} \frac{\|x\|}{\|y\|} \geq \frac{2}{\pi} \tan^{-1} \frac{1}{1 - \gamma}.$$

In order to bound to the latter quantity, we consider that $\tan^{-1}(1 - \gamma)^{-1}$ is a convex function, and its derivative in $\gamma = 0$ is $\frac{1}{2}$. So we obtain

$$\tan^{-1} \frac{1}{1 - \gamma} \geq \tan^{-1} 1 + \frac{1}{2}\gamma = \frac{\pi}{4} + \frac{1}{2}\gamma$$

from which the lemma follows. \square

Intuitively, we want to fix a set V of directions, and compare the lengths of vectors by comparing the lengths of their respective projections over elements of V , and making a majority vote. For two vectors x, y we will write that

$$x \triangleright_V y \Leftrightarrow |V \cap \mathcal{L}_{x,y}| \geq \frac{1}{2}|V|. \tag{1}$$

This means that the projection of x is longer than the projection of y over at least half of the vectors in V . Note that $x \triangleright_V y$ and $y \triangleright_V x$ can simultaneously hold with respect to the same set V .

For a fixed set of directions V , consider two arbitrary vectors x and y such that x is significantly longer than y . Lemma 1 says that the set $\mathcal{L}_{x,y}$ is large, so we can hope that in V there are enough vectors of $\mathcal{L}_{x,y}$ so that $x \triangleright_V y$. However, we want this to hold for *any* vectors x, y , with respect to the *same* (fixed) set V . For this reason we cannot use directly Lemma 1, but we will use some VC-dimension techniques. For a detailed treatment of this topic see [35, 36, 1, 29], or the paper which introduced the use of VC-dimension in computational geometry [22]. Here, we will present only the definitions needed.

A *range space* is a pair $(\mathcal{P}, \mathcal{R})$, where $\mathcal{P} = (\Omega, \mathcal{F}, \mu)$ is a probability space, and $\mathcal{R} \subseteq \mathcal{F}$ is a collection of measurable (w.r.t. the measure μ) subsets of Ω . A finite set $A \subseteq \Omega$ is said to be *shattered* by \mathcal{R} if every subset of A can be expressed in the form $A \cap R$ for some $R \in \mathcal{R}$. The *VC-dimension* $VC\text{-dim}(\mathcal{R})$ of the range space $(\mathcal{P}, \mathcal{R})$ is the maximum size of a set that can be shattered by \mathcal{R} (hence no set of size $> VC\text{-dim}(\mathcal{R})$ can be shattered).

There is a natural identification between collections of sets and families of binary-valued function: to each subset $R \in \mathcal{R}$ corresponds its indicator function $f_R(x) : \Omega \rightarrow \{0, 1\}$ such that $f_R(x) = 1 \Leftrightarrow x \in R$. This identification is useful to express combinations of range spaces in the following manner.

For an integer $k \geq 2$, let $u : \{0, 1\}^k \rightarrow \{0, 1\}$ and $f_1, \dots, f_k : \mathcal{R} \rightarrow \{0, 1\}$ be binary-valued functions (u represents the combination operator and the f_i 's are indicator

functions). Define $u(f_1, \dots, f_k) : \mathcal{R} \rightarrow \{0, 1\}$ as the binary-valued function $x \mapsto u(f_1(x), \dots, f_k(x))$. Finally, if $\mathcal{A}_1, \dots, \mathcal{A}_k$ are families of binary-valued functions, we define $\mathcal{U}(\mathcal{A}_1, \dots, \mathcal{A}_k)$ to be the family of binary-valued functions $\{u(f_1, \dots, f_k) : f_i \in \mathcal{A}_i \ \forall i\}$. In this way we can obtain $\mathcal{A} \oplus \mathcal{B} = \{A \cup B : A \in \mathcal{A}, B \in \mathcal{B}\}$ by choosing $u(a, b) = \max\{a, b\}$, and $\mathcal{A} \odot \mathcal{B} = \{A \cap B : A \in \mathcal{A}, B \in \mathcal{B}\}$ by choosing $u(a, b) = a \cdot b$.

We will use the following theorem by Vidyasagar [36, Theorem 4.3], which improves on a result by Dudley [10, 11]:

Theorem 1. *If $VC\text{-dim}(\mathcal{A}_i)$ is finite for each i , then $\mathcal{U} = \mathcal{U}(\mathcal{A}_1, \dots, \mathcal{A}_k)$ also has finite VC-dimension, and $VC\text{-dim}(\mathcal{U}) < \alpha_k \cdot d$, where $d = \max_i VC\text{-dim}(\mathcal{A}_i)$, and α_k is the smallest integer such that $k < \alpha_k / \log_2(e\alpha_k)$.*

Some values of α_k are the following:

k	2	3	4	5	6	7	8	9	10
α_k	10	17	25	33	41	50	59	68	78

We consider the probability space \mathcal{P} made by the set of directions S^{d-1} , with the uniform distribution, where $\mu(X) = |X|/|S^{d-1}|$; we take as \mathcal{R} the collection of all the sets $\mathcal{L}_{x,y}$. The following lemma bounds the VC-dimension of this range space:

Lemma 2. *The VC-dimension of the range space $(\mathcal{P}, \mathcal{R})$ defined above is strictly less than $25(d + 1)$.*

Proof. For a vector $u \in \mathbb{R}^d$, let H_u denote the closed hemisphere $\{v \in S^{d-1} : u \cdot v \geq 0\}$, and let \mathcal{H} be the collection of all such hemispheres. Then any set $\mathcal{L}_{x,y}$ can be expressed as the Boolean combination $(H_{x-y} \cap H_{x+y}) \cup (H_{-x-y} \cap H_{-x+y})$. Thus we have $\mathcal{R} \subseteq \mathcal{U}$, where $\mathcal{U} = (\mathcal{H} \odot \mathcal{H}) \oplus (\mathcal{H} \odot \mathcal{H})$. A theorem by Radon [4] implies that $VC\text{-dim}(\mathcal{H}) \leq d + 1$. We then apply Theorem 1, with $u(a, b, c, d) = \max\{a \cdot b, c \cdot d\}$, and the fact that $\alpha_4 = 25$ to obtain that $VC\text{-dim}(\mathcal{R}) \leq VC\text{-dim}(\mathcal{U}) < 25(d + 1)$. \square

A subset A of Ω is called a γ -approximation (or γ -sample) for a range space $(\mathcal{P}, \mathcal{R})$ is for every $R \in \mathcal{R}$ the following relation holds:

$$\left| \frac{|R \cap A|}{|A|} - \mu(R) \right| \leq \gamma. \tag{2}$$

This means that A can be used to obtain a good estimate of the measure of any set $R \in \mathcal{R}$. The main result, which follows from Lemma 2 and from a fundamental theorem by Vapnik and Chervonenkis, is the following.

Lemma 3. *With probability at least $1 - \delta$ a set V of cardinality*

$$f(\gamma, \delta) = \frac{16}{\gamma^2} \left(25(d + 1) \log \frac{400(d + 1)}{\gamma^2} + \log \frac{4}{\delta} \right) = \Theta(d \log d)$$

of vectors chosen uniformly at random from S^{d-1} is a γ -approximation for the range space defined above.

The lemma above permits to make comparisons in the following way: choose a set V of $f(\varepsilon/\pi, \delta)$ random directions. With probability $1 - \delta$ it is a (ε/π) -approximation. For any x, y with $(1 - \varepsilon)\|x\| \geq \|y\|$, we have that $\mu(\mathcal{L}_{x,y}) \geq \frac{1}{2} + \varepsilon/\pi$, by Lemma 1. By definition of γ -approximation we have that $|\mathcal{L}_{x,y} \cap V| \geq \frac{1}{2}|V|$, so $x \triangleright_V y$ by definition (1). This is the idea which we are going to use in the following section. However, in order to save computations, we will make comparisons using random subsets of the fixed set V .

3. An algorithm for farthest-point queries

We will present in this section a (ε, δ) -approximation scheme for computing the farthest site of a query point q . This means that with probability $1 - \delta$ the algorithm gives an answer which is within a factor $1 - \varepsilon$ from the optimal one. The parameters ε and δ are chosen by the user before the algorithm starts.

Let $P = \{p_1, \dots, p_n\}$ the set of given sites, and $q \in \mathbb{R}^d$ the query point. For simplicity we will assume that n is a power of 2. Let p^* be the farthest-site from q , i.e., the point such that

$$d(p^*, q) = \max_{p_i \in P} d(p_i, q),$$

where $d(p, q) = \|p - q\|$ is the standard Euclidean metric in \mathbb{R}^d . Let Z_ε be the set of sites that are far from q ‘almost’ like p^* :

$$Z_\varepsilon = \{p_i \in P : d(p_i, q) \geq (1 - \varepsilon)d(p^*, q)\}.$$

The purpose of the algorithm is to give, with probability at least $1 - \delta$, an element of Z_ε .

3.1. Building the data structure

The preprocessing stage is the following. Let $\varepsilon_0 = \log(1 + \varepsilon)/\log n$. We choose randomly L vectors from S^{d-1} , where $L = f(\varepsilon_0/12, \delta) = \Theta(d \log d \log^2 n \log \log n)$. These vectors can be obtained for example by the method described by Knuth [26, p. 130]. Let $V = \{v_1, \dots, v_L\}$ the set of L directions generated above. The data structure is a matrix M , of dimension $L \times n$, where $M[i, j] = v_i \cdot p_j$.

3.2. Processing a query

We first define the following relation between sites of P : for a set of directions $\Gamma \subseteq V$, we say that $p_i \triangleright_\Gamma p_j$ if the projection of p_i is farther than the projection of p_j

from the projection of q , for at least half of the directions in Γ . Formally:

$$p_i \succ_{\Gamma} p_j \Leftrightarrow (p_i - q) \triangleright_{\Gamma} (p_j - q).$$

We will use this relation to build a ‘tournament’ between sites, by making comparisons with respect to a fixed set Γ : if $p_i \succ_{\Gamma} p_j$ and $p_j \not\succeq_{\Gamma} p_i$, the winner of the comparison is p_i ; if $p_i \not\succeq_{\Gamma} p_j$ and $p_j \succ_{\Gamma} p_i$, the winner is p_j ; finally, if both $p_i \succ_{\Gamma} p_j$ and $p_j \succ_{\Gamma} p_i$, the winner is chosen arbitrarily. Note that the above description is complete, i.e., it cannot occur that $p_i \not\succeq_{\Gamma} p_j$ and $p_j \not\succeq_{\Gamma} p_i$.

Phase A. We first extract a random subset $\Gamma \subseteq V$ of cardinality $c_1 \log^3 n$ (the value of c_1 will be given in Lemma 7). We make extractions with replacement, so we permit Γ to be a multi-set. Let $b = \log|\Gamma| = \Theta(\log \log n)$. We assume for simplicity that b has an integer value.

We compute the values $v \cdot q$ for all $v \in \Gamma$. Then we build a complete binary tree T of depth $\log n$. We associate randomly the sites in P to the leaves of T . To every internal node x at height $1 \leq h \leq b$, we associate a random subset $\Gamma_x \subseteq \Gamma$ of size $c'_2 + c_2 h$ (the appropriate values for c_2 and c'_2 will be given in Lemma 8). To higher nodes, with height $h > b$, we associate the entire set Γ .

Now we make a tournament between the sites, proceeding from the leaves towards the root of T : to each internal node x we associate the winner of the comparison between the sites associated to its children, with respect to the set Γ_x . Let \tilde{p}_A be the winner of this tournament, i.e., the site which at the end of Phase A will be associated to the root of T .

Phase B. Independently from all the above, we randomly choose a subset $P_0 \subseteq P$ of size $c_3 \log^3 n$, where c_3 will be defined in Lemma 6. We compute the distances between q and the sites in P_0 . Let \tilde{p}_B be the site of P_0 farthest from q .

The algorithm finishes returning the site among \tilde{p}_A and \tilde{p}_B that is the farthest from q .

3.3. Correctness

We now prove that the algorithm correctly computes an element of Z_{ε} . First of all, observe that, as a consequence of Lemma 3, the set V created during the preprocessing is an $(\varepsilon_0/12)$ -approximation of the range space $(\mathcal{P}, \mathcal{R})$, with probability $1 - \delta$. In the following, we will assume that this event occurred.

Lemma 4. *Let V be an $(\varepsilon_0/12)$ -approximation, and v a random element of V . Let $q \in \mathbb{R}^d$, $p_i, p_j \in P$, and suppose that $(1 - \gamma)d(p_i, q) \geq d(p_j, q)$, with $\varepsilon_0 \leq \gamma \leq \frac{1}{2}$. The probability that $|v \cdot (p_i - q)| \geq |v \cdot (p_j - q)|$ is at least $\frac{1}{2} + \gamma/5$.*

Proof. Let $x = p_i - q$ and $y = p_j - q$. The measure $z_{x,y}$ of $\mathcal{L}_{x,y}$ is at least $\frac{1}{2} + \gamma/\pi$, by Lemma 1. As V is an $(\varepsilon_0/12)$ -approximation, the probability that an element of V belongs to $\mathcal{L}_{x,y}$ is at most $\varepsilon_0/12$ less than that value. So this probability is at least $z_{x,y} - \varepsilon_0/12 \geq \frac{1}{2} + \gamma/5$, and the lemma follows. \square

Lemma 5. *Suppose $(1 - \gamma)d(p_i, q) \geq d(p_j, q)$, and let Γ' be a random subset of Γ , of cardinality k . The probability that $p_j \succ_{\Gamma'} p_i$ is at most $e^{-\gamma^2 k/36}$.*

Proof. Let $x = p_i - q$ and $y = p_j - q$. The relation $p_j \succ_{\Gamma'} p_i$ is equivalent to $y \triangleright_{\Gamma'} x$, that means $|\Gamma' \cap \mathcal{L}_{y;x}| \geq \frac{1}{2}k$. As $S^{d-1} \setminus \mathcal{L}_{y;x} \subseteq \mathcal{L}_{x;y}$, we need only to bound the probability that $|\Gamma' \cap \mathcal{L}_{x;y}| \leq \frac{1}{2}k$. Define k Bernoulli random variables X_1, \dots, X_k , where X_r is 1 if the r -th vector in Γ' belongs to $\mathcal{L}_{x;y}$, and 0 otherwise. By application of Lemma 4 we obtain that the probability of success is $\text{PR}[X_r = 1] \geq \frac{1}{2} + \gamma/5$. Defining the random variable $X = \sum_r X_r = |\Gamma' \cap \mathcal{L}_{x;y}|$, the event we are interested in is $X \leq \frac{1}{2}k$. The lemma follows by some calculations and by application of the Chernoff bound (see e.g. [28]).

Here are the details. The expected value of X is $E[X] \geq k(\frac{1}{2} + \gamma/5)$. Using the fact that $\frac{1}{2} < (1 - \gamma/3)(\frac{1}{2} + \gamma/5)$ we have:

$$\text{PR} \left[X \leq \frac{1}{2}k \right] \leq \text{PR} \left[X < \left(1 - \frac{\gamma}{3}\right) \left(\frac{1}{2} + \frac{\gamma}{5}\right) k \right] = \text{PR} \left[X < \left(1 - \frac{\gamma}{3}\right) E[X] \right].$$

The Chernoff bound (see e.g. [28, Theorem 4.2]) affirms that $\text{PR}[X < (1 - \delta)E[X]] < e^{-\delta^2 E[X]/2}$. Applying this result we obtain

$$\text{PR} \left[X < \left(1 - \frac{\gamma}{3}\right) E[X] \right] < \exp \left(- \left(\frac{\gamma}{3}\right)^2 k \left(\frac{1}{2} + \frac{\gamma}{5}\right) \frac{1}{2} \right) \leq e^{-\gamma^2 k/36},$$

where we also used that $e^{-p} \leq e^{-1/2}$ for $p \geq \frac{1}{2}$. \square

The motivation for Phase B is that if in the set Z_e there are too many sites, one of them could eliminate p^* during the early stage of the tournament, and then be eliminated by some site not in Z_e . On the other hand, if Z_e is sufficiently small, there is a high probability that p^* will reach at least level b of the tournament tree: in this case, we can show that the winner \tilde{p}_A is an element of Z_e . Intuitively, at each comparison the winner can be slightly closer to q than the loser: if p^* is eliminated too early in the competition, all these small errors can take us too close to q ; if instead p^* reaches at least level b , the final error is small with high probability.

We begin by proving that if Z_e is large, then Phase B succeeds in finding an approximate farthest-point of q with high probability:

Lemma 6. *If $|Z_e| \geq \gamma_1 n / \log^3 n$, then $\tilde{p}_B \in Z_e$, with probability at least $1 - \delta$.*

Proof. When we randomly select the elements of P_0 from P , the probability of taking an element not in Z_e is less than $1 - \gamma_1 / \log^3 n$. The probability that no element of P_0 belongs to Z_e is less than that quantity raised to power $|P_0| = c_3 \log^3 n$. We can choose $c_3 \equiv c_3(\gamma_1, \delta, n)$ such that $(1 - \gamma_1 / \log^3 n)^{c_3 \log^3 n} \leq \delta$. If we choose such a c_3 , with probability at least $1 - \delta$ there is at least one element of Z_e in P_0 , so \tilde{p}_B certainly belongs to Z_e . \square

Let us see now what happens if Z_ε is small. Define the constant $\gamma_2 \equiv \gamma_2(\varepsilon)$ such that $e^{-\gamma_2} \geq 1 - \varepsilon$. We denote by \mathcal{E}_1 the following event:

$$\mathcal{E}_1 \equiv \left\{ \exists p_i, p_j \in P, \quad i \neq j: p_j \succ_{\Gamma} p_i, \quad \left(1 - \frac{\gamma_2}{\log n} \right) d(p_i, q) \geq d(p_j, q) \right\},$$

that is, the event that there exist two sites, whose distances from q differ significantly, for which the comparison based on projections on the vectors of Γ gives the wrong answer. We can give a bound to the probability that this event occurs:

Lemma 7. *The probability of event \mathcal{E}_1 is less or equal to $\delta/3$.*

Proof. There are $n(n - 1) \leq n^2$ ordered pairs of sites to consider: for each pair we apply Lemma 5, with $\gamma = \gamma_2 / \log n$ and $k = |\Gamma| = c_1 \log^3 n$, to bound the probability of error. The probability of \mathcal{E}_1 is bounded by the sum of these probabilities. The lemma follows by defining suitably the value of $c_1 \equiv c_1(\gamma_2, \delta, n)$ in such a way that $n^2 \cdot e^{-c_1 \gamma_2^2 \log n / 36} \leq \delta/3$. \square

Let us denote by \mathcal{E}_2 the event “ p^* does not reach level b in the tournament tree”, that is, it is not assigned, during the tournament in Phase A, to any node at level b .

Lemma 8. *If $|Z_\varepsilon| < \gamma_1 n / \log^3 n$, the probability of \mathcal{E}_2 is less or equal to $2\delta/3$.*

Proof. Consider the leaf of T to which p^* was assigned, the node x at level b along the path from this leaf to the root (x is the node to which p^* should be assigned, if \mathcal{E}_2 does not occur), and the sub-tree T^x rooted at x .

We first of all exclude the possibility that in T^x there be sites of Z_ε other than p^* . The number of leaves in T^x is $2^b = |\Gamma| = c_1 \log^3 n$, so the probability that $|T^x \cap Z_\varepsilon| = 1$ is

$$\binom{n - |Z_\varepsilon|}{2^b - 1} \binom{n}{2^b - 1}^{-1}.$$

By defining suitably $\gamma_1 \equiv \gamma_1(c_1, \delta, n)$ we can make this probability greater than $1 - \delta/3$. So we can assume that p^* is the only site of Z_ε in the sub-tree T^x . We now show that in this case it reaches level b with probability $1 - \delta/3$.

If p^* reached level $h - 1$, it is compared, at level h , with an element not in Z_ε , made using $|I_x| = c'_2 + c_2 h$ vectors of Γ . By application of Lemma 5, with $\gamma = \varepsilon$, the probability that p^* loses the comparison is bounded by

$$e^{-\varepsilon^2(c'_2 + c_2 h) / 36} = e^{-\varepsilon^2 c'_2 / 36} (e^{-\varepsilon^2 c_2 / 36})^h.$$

We can define the constants $c'_2 \equiv c'_2(\varepsilon, \delta)$ and $c_2 = c_2(\varepsilon)$ such that the left factor is less than $\delta/3$, and the right factor is less than 2^{-h} .

Summing these quantities for $h = 1, \dots, b$ we finally obtain that the probability that p^* is eliminated is bounded by $\delta/3$. Together with the fact that the probability that more than one site of Z_ε belongs to T^x is less than $\delta/3$, this implies the lemma. \square

We can now prove that if Z_ε is small, Phase A succeeds in finding a good answer:

Lemma 9. *If $|Z_\varepsilon| < \gamma_1 n / \log^3 n$, then $\tilde{p}_A \in Z_\varepsilon$, with probability at least $1 - \delta$.*

Proof. With probability $1 - \delta$ neither \mathcal{E}_1 nor \mathcal{E}_2 occur. In this case we show by induction on $h \geq b$ that there exists a site $p_{(h)}$, assigned to a node at level h , such that

$$d(p_{(h)}, q) \geq \left(1 - \frac{\gamma_2}{\log n}\right)^{h-b} d(p^*, q) \tag{3}$$

and this implies that $\tilde{p}_A \in Z_\varepsilon$.

The fact that \mathcal{E}_2 does not occur guarantees that p^* reached level b , so we can start the induction, for $h = b$, with $p_{(h)} = p^*$. Suppose now that (3) holds at level h , and let $p'_{(h)}$ be the adversary of $p_{(h)}$. If $d(p'_{(h)}, q) \geq (1 - \gamma_2 / \log n) d(p_{(h)}, q)$, then both $p_{(h)}$ and $p'_{(h)}$ satisfy (3) for level $h + 1$, so relation (3) holds in any case. If instead $d(p'_{(h)}, q) < (1 - \gamma_2 / \log n) d(p_{(h)}, q)$, as we are assuming that \mathcal{E}_1 does not occur, certainly $p_{(h)}$ wins the comparison, so $p_{(h+1)} = p_{(h)}$ and relation (3) is valid for level $h + 1$.

From the inductive argument we deduce that (3) holds for $h = \log n$, where $p_{(\log n)} = \tilde{p}_A$ is the site assigned to the root of T . Being $b \geq 1$, from the fact that $(1 - x/m)^{m-1} \geq e^{-x}$, and from the definition of γ_2 , we obtain

$$d(\tilde{p}_A, q) \geq \left(1 - \frac{\gamma_2}{\log n}\right)^{\log n - 1} d(p^*, q) \geq e^{-\gamma_2} d(p^*, q) \geq (1 - \varepsilon) d(p^*, q),$$

that is what we wanted to prove. \square

The correctness of the algorithm follows from both Lemmas 6 and 9.

3.4. Complexity

The preprocessing consists in the computation of a distance, which requires time $O(d)$, for each element of the matrix M . The time needed is thus $O(dLn) = O(d^2 n \log d \log^2 n)$. The space required to store M is $O(Ln) = O(dn \log d \log^2 n)$.

To answer a query, we first compute $v \cdot q$ for $v \in \Gamma$. The time required for this is $O(|\Gamma| \cdot d) = O(d \log^3 n)$. Each comparison of type $p_i \geq_{\Gamma_x} p_j$ requires time $O(|\Gamma_x|)$, because all values $v \cdot p_i$ are already stored in M . The number of nodes at level h is $n2^{-h}$, so the total number of operations at levels $1 \leq h \leq b$ is

$$\sum_{h=0}^b (c'_2 + c_2 h) n 2^{-h} = c'_2 n \sum 2^{-h} + c_2 n \sum h 2^{-h} \leq 2(c'_2 + c_2) n = O(n).$$

For levels $h > b$ each comparison is made using the entire set Γ , where $|\Gamma| = 2^b$. We have

$$\sum_{h=b+1}^{\log n} |\Gamma| \cdot n 2^{-h} = \sum_{h=b+1}^{\log n} |\Gamma| \cdot n 2^{b-h} 2^{-b} = n \sum_{h=b+1}^{\log n} 2^{b-h} \leq n.$$

We thus obtain that the total cost of Phase A is $O(n + d \log^3 n)$. Phase B requires time $O(d \log^3 n)$. We conclude that the time required to answer a query is $O(n + d \log^3 n)$. We summarize the properties of the farthest-point algorithm in the following theorem:

Theorem 2. *With probability $1 - \delta$ the set V created in the preprocessing is an $(\varepsilon_0/12)$ -approximation. In this case, the probability that the algorithm returns an element of Z_ε is, for each query, at least $1 - \delta$. The complexity of the algorithm is $O(d^2 \log d \cdot n \log^2 n)$ for the preprocessing and $O(n + d \log^3 n)$ for answering a query.*

4. Computing the diameter of a point set

We now show how to use the algorithm for farthest-point queries to find the diameter d_P of a set of points $P = \{p_1, \dots, p_n\}$ in \mathbb{R}^d . First of all, we use a dimension-reduction technique, by projecting all the points on a random subspace of dimension $k = O(\varepsilon^{-2} \log n)$. Let P' the resulting set of points in \mathbf{R}^k and $d_{P'}$ its diameter. The JL Lemma (see [24, 16]) affirms that with high probability all inter-point distances are preserved, so that $(1 + \varepsilon/2)d_P \geq d_{P'} \geq (1 - \varepsilon/2)d_P$.

Next, we compute the set V and the matrix M as in the preprocessing stage of the farthest-point algorithm, using $\varepsilon/2$ as approximation parameter. For each $p'_i \in P'$, we perform a farthest-point query where $q = p'_i$ and the sites are the remaining points in P' ; let $F_\varepsilon(p'_i)$ the point returned by the algorithm described in Section 3. We compute all the distances $d(p'_i, F_\varepsilon(p'_i))$. Let $\tilde{d}_{P'}$ the maximum distance computed in this way. Our main result is the following:

Theorem 3. *Let d_P be the diameter of the point set P , and $\tilde{d}_{P'}$ the value computed by the algorithm above. Then, with probability $1 - \delta$, it holds that $(1 + \varepsilon/2)d_P \geq \tilde{d}_{P'} \geq (1 - \varepsilon)d_P$. The complexity of the algorithm is $O(nd \log n + n^2)$.*

Proof. Let (p^*, q^*) be a pair of points in P' such that $d(p^*, q^*) = d_{P'}$. When we make the farthest-point query with $q = p^*$, the point $F_\varepsilon(p^*)$ is, with probability $1 - \delta$, such that $d(p^*, F_\varepsilon(p^*)) \geq (1 - \varepsilon/2)d_{P'}$. The value $\tilde{d}_{P'}$ computed in the last step is $\tilde{d}_{P'} \geq (1 - \varepsilon/2)d_{P'} \geq (1 - \varepsilon/2)^2 d_P \geq (1 - \varepsilon)d_P$. Moreover, the distances computed by this algorithm are certainly less than $d_{P'}$.

The time to perform the projection at the first step is $O(nd \log n)$. Then we perform the preprocessing and n queries of the algorithm in Section 3, in dimension $O(\log n)$. As a consequence of Theorem 2, we obtain that the overall complexity of the algorithm is $O(nd \log n + n^2)$. \square

Note that by taking the points whose projected distance is $\tilde{d}_{P'}$, and computing their distance in the original d -dimensional space, we can obtain a value $\tilde{d}'_{P'}$ such that $d_P \geq \tilde{d}'_{P'} \geq (1 - \varepsilon)d_P$.

5. Conclusions

We have shown how the application of randomized projection techniques results in an algorithm which computes an approximation of the diameter in time $O(n^2)$, independent of the dimension, for dimensions up to $n/\log n$. The algorithm is simple (although its analysis is not) and has a potential for an impact on practical application.

Acknowledgements

We thank Piotr Indyk and Jon Kleinberg for pointing out the paper [5], and Yuval Rabani for sending us a preliminary version of it. We thank Jirka Matoušek for his comments and suggestions and Andrew Yao for clarifications on his algorithms in [37]. We thank an anonymous referee who kindly pointed out a mistake in a previous version of this paper.

References

- [1] N. Alon, J.H. Spencer, *The Probabilistic Method*, Wiley, New York, 1992.
- [2] D.M. Avis, A survey of heuristics for the weighted matching problem, *Networks* 13 (1983) 475–493.
- [3] S.N. Bespamyatnikh, An efficient algorithm for the three-dimensional diameter problem, in: *Proc. 9th Symp. on Discrete Algorithms*, 1998, pp. 137–146.
- [4] A. Blumer, A. Ehrenfeucht, D. Haussler, M.K. Warmuth, Learnability and the Vapnik–Chervonenkis dimension, *J. ACM* 4 (36) (1989) 929–965.
- [5] A. Borodin, R. Ostrovsky, Y. Rabani, Subquadratic approximation algorithms for clustering problems in high dimensional spaces, in: *Proc. 31st ACM Symp. on the Theory of Computing*, 1999, pp. 435–444.
- [6] A.Z. Broder, S.C. Glassman, M.S. Manasse, G. Zweig, Syntactic clustering of the Web, *Computer Networks and ISDN Systems* 29 (8–13) (1997) 1157–1166.
- [7] B. Chazelle, H. Edelsbrunner, L.J. Guibas, M. Sharir, Diameter, width, closest line pair, and parametric searching, *Discrete and Comput. Geom.* 10 (1993) 183–196.
- [8] K.L. Clarkson, P.W. Shor, Applications of random sampling in computational geometry, II, *Discrete and Comput. Geom.* 4 (1989) 387–421.
- [9] D. Coppersmith, S. Winograd, Matrix multiplication via arithmetic progressions, in: *Proc. 19th ACM Symp. on Theory of Computing*, 1987, pp. 1–6.
- [10] R.M. Dudley, Central limit theorems for empirical measures, *Ann. Probab.* 6 (1978) 899–929.
- [11] R.M. Dudley, A course on empirical processes, in: *Lecture Notes in Mathematics*, vol. 1097, Springer, Berlin, 1984, pp. 1–142.
- [12] O. Egecioğlu, B. Kalantari, Approximating the diameter of a set of points in the Euclidean space, *Inform. Process. Lett.* 32 (1989) 205–211.
- [13] P. Erdős, On sets of distances of n points, *Amer. Math. Monthly* 53 (1946) 248–250.
- [14] P. Erdős, On sets of distances of points in Euclidean spaces, *Magyar Tud. Akad. Mat. Kutató Int. Közl.* 5 (1960) 165–169.
- [15] M. Flickner, H. Sawhney, W. Niblack, J. Ashley, Q. Huang, B. Dom, M. Gorkani, J. Hafner, D. Lee, D. Petkovic, D. Steele, P. Yanker, Query by image and video content: The QBIC system, *Computer* 28 (9) (1995) 23–32.
- [16] P. Frankl, H. Maehara, The Johnson–Lindenstrauss lemma and the sphericity of some graphs, *J. Combin. Theory Series B* 44 (1998) 355–362.
- [17] P. Gritzmann, V. Klee, Inner and outer j -radii of convex bodies in finite dimensional normed spaces, *Discrete and Comput. Geom.* 7 (1992) 255–280.

- [18] P. Gritzmann, V. Klee, Computational complexity of inner and outer j -radii of polytopes in finite-dimensional normed spaces, *Math. Program* 59 (1993) 163–213.
- [19] P. Gritzmann, V. Klee, On the complexity of some basic problems in computational convexity: I. Containment problems, *Discrete Math.* 136 (1994) 129–174.
- [20] V.N. Gudivada, V.V. Raghavan, Guest editors' introduction: Content-based image retrieval systems, *Computer* 28 (9) (1995) 18–22.
- [21] J.A. Hartigan, *Clustering Algorithms*, Wiley, New York, 1975.
- [22] D. Haussler, E. Welzl, ε -nets and simplex range queries, *Discrete and Comput. Geom.* 2 (1987) 127–151.
- [23] C.E. Jacobs, A. Finkelstein, D.H. Salesin, Fast multiresolution image querying, in: R. Cook (Ed.), *SIGGRAPH 95 Conference Proc.*, Addison-Wesley, Reading, MA, 1995, pp. 277–286.
- [24] W.B. Johnson, J. Lindenstrauss, Extensions of Lipschitz mappings into a Hilbert space, *Contemporary Mathematics* 26 (1984) 189–206.
- [25] J.M. Kleinberg, Two algorithms for nearest-neighbor search in high dimensions, in: *Proc. 29th ACM Symp. on the Theory of Computing*, 1997, pp. 599–608.
- [26] D. Knuth, *The Art of Computer Programming*, Addison-Wesley, Reading, MA, 1973.
- [27] J. Matoušek, O. Schwarzkopf, A deterministic algorithm for the three-dimensional diameter problem, in: *Proc. 25th ACM Symp. on the Theory of Computing*, 1993, pp. 478–484.
- [28] R. Motwani, P. Raghavan, *Randomized Algorithms*, Cambridge University Press, Cambridge, 1995.
- [29] K. Mulmuley, *Computational Geometry: An Introduction Through Randomized Algorithms*, Prentice-Hall, Englewood Cliffs, NJ, 1994.
- [30] J. O'Rourke, G. Toussaint, Pattern recognition, in: J.E. Goodman, J. O'Rourke (Eds.), *Handbook of Discrete and Computational Geometry*, CRC Press, Boca Raton, FL, 1997, pp. 797–813.
- [31] F.P. Preparata, M.I. Shamos, *Computational Geometry: an Introduction*, Springer, New York, 1985.
- [32] E.A. Ramos, Construction of 1-d lower envelopes and applications, in: *Proc. of the 13th Symp. on Computational Geometry*, 1997, pp. 57–66.
- [33] S. Straszewicz, Sur un problème géométrique de P. Erdős, *Bull. Acad. Polon. Sci. Cl. III* 5 (1957) 39–40.
- [34] K.J. Supowit, New techniques for some dynamic closest-point and farthest-point problems, in: *Proc. 1st ACM-SIAM Symp. on Discrete Algorithms*, 1990, pp. 84–90.
- [35] V.N. Vapnik, A.Y. Chervonenkis, On the uniform convergence of relative frequencies of events to their probabilities, *Theory Probab. Appl.* 16 (1971) 264–280.
- [36] M. Vidyasagar, *A Theory of Learning and Generalization: with Applications to Neural Networks and Control Systems*, Springer, London, 1997.
- [37] A.C. Yao. Personal Communication, Also in "On computing the distance matrix of n points in k dimensions", TR IBM San Jose Research Center, 1982.
- [38] A.C. Yao, On constructing minimum spanning trees in k -dimensional spaces and related problems, *SIAM J. Comput.* 11 (4) (1982) 721–736.