

Approximate Nearest Neighbor Queries among Parallel Segments

Ioannis Z. Emiris*

Theocharis Malamatos†

Elias Tsigaridas‡

Abstract

We develop a data structure for answering efficiently approximate nearest neighbor queries over a set of parallel segments in three dimensions. We connect this problem to approximate nearest neighbor searching under weight constraints and approximate nearest neighbor searching on historical data in any dimension and we give efficient solutions for these as well.

1 Introduction

Nearest neighbor searching is a fundamental geometric problem with applications in many areas. For high dimensions there are no known efficient exact solutions and thus approximate solutions to the problem have been studied. Let $d(p, q)$ denote the euclidean distance between points p, q . Given a set P of points in \mathbb{R}^d and a parameter $\varepsilon > 0$, we say that a point p of P is an ε -approximate nearest neighbor (ε -NN) to a point q if $d(p, q) \leq (1 + \varepsilon)d(q', q)$ where q' is a nearest point to q in P . Arya et al. [3] have shown how to find efficiently an ε -NN to any given query point in constant dimensions and Indyk and Motwani [6] presented efficient methods for high dimensions. See [5] for more references.

An interesting generalization of the problem arises if we replace the point set P with a set of objects O . For this version there are only few results known. When O is a set of disjoint polyhedra in three dimensions Koltun and Sharir [7] presented a data structure of near quadratic size that can answer an ε -NN query in $O(\log(n/\varepsilon))$ time. In three dimensions again, Wang [9] showed how to answer ε -NN queries when O is a set of triangles, segments, and points in a convex position in $O(\log^2 n/\varepsilon^2)$ query time and using $O(n/\varepsilon^2)$ space. In high dimensions, Magen [8] provided an algorithm for a set of k -flats with query time polynomial in $d, \log n$ and $1/\varepsilon$ but non-polynomial space.

In this paper we consider the case where O is a set of parallel segments. We give two solutions for the

problem. The second solution enhances the first using results in Sec. 3. That section is motivated by the *cheap gas station* problem. Given n gas stations (sites) and a car at a position q (query point) we want to find a gas station that is closest or approximately closest to q (since exact distance is not so important) which sells gas for at most w euros. In Sec. 4 we combine the results of the previous sections and present a data structure for answering time-dependent ε -NN queries in any fixed dimension.

2 Methods for parallel segments

Let S be a set of n disjoint parallel segments in \mathbb{R}^3 . We assume w.l.o.g. that all segments in S are parallel to the x -axis. Let P be the set of the $2n$ endpoints of the segments in S . Let q be a query point in \mathbb{R}^3 , s be a segment of S nearest to q , and p be the point of s nearest to q . Observe that p is either one of the endpoints of s or that the segment qp is perpendicular to s . Let H_q be the plane passing through q that is parallel to the yz plane. Note that if p is interior to s then p is one of the points in $H_q \cap S$. (See Fig. 1.)

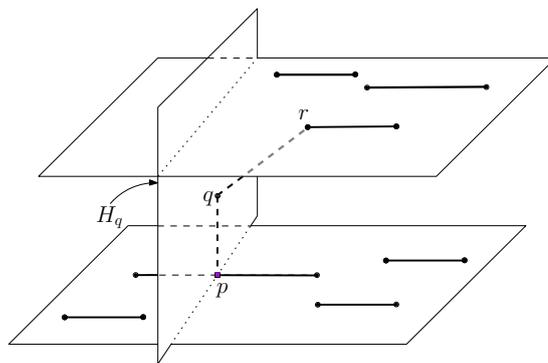


Figure 1: Parallel segments lying in two parallel planes in \mathbb{R}^3 and a query point q . The closest endpoint to q is r but the closest point to q is p .

It follows that to find an ε -NN to q in S it suffices to (a) find an ε -NN to q in P , (b) find an ε -NN to q in the set $H_q \cap S$ and then report whichever of the two is nearest to q . For solving (a) we use an (t, ε) -approximate Voronoi diagram (AVD) on set P with $t = O(1/\varepsilon)$ together with the associated data structure [2]. This structure has $O(n)$ space and it returns an ε -NN to any q in $O(\log n + 1/\varepsilon)$ time. For solving (b) we present two methods which are both based on

*Dept. of Informatics and Telecoms, National and Kapodistrian University of Athens, Greece.

†Dept. of Computer Science and Technology, University of Peloponnese, Greece.

‡Dept. of Computer Science, Aarhus University, Denmark and Dept. of Computer Science and Technology, University of Peloponnese, Greece. Partially supported by a research grant from the Danish Council of Independent Research.

a variation of the well-known interval tree [4]. The second method is more complicated but leads to a significant improvement over the first.

2.1 First method

The construction of our interval tree T on S proceeds as follows. Let x_m be the median among all x -coordinates of P . Let H_m be the plane passing through point $(x_m, 0, 0)$ and which is parallel to the yz plane. We store x_m at the root of T . We partition the set of segments S into three sets S_ℓ , S_m , and S_r where each set consists of the segments that lie on the left of H_m , that intersect H_m , and that lie on the right of H_m , respectively. We continue the construction of the tree T recursively on S_ℓ and S_r . (If S_ℓ or S_r is the empty set clearly we get a leaf.) The two roots of the trees built on S_ℓ and on S_r become the left and right child of the root of T , respectively. See Fig. 2 for an example. For the set S_m we build an

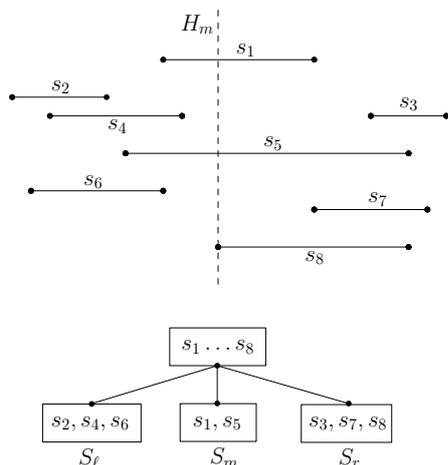


Figure 2: A projection of a set of segments in xz plane. The left endpoint of s_8 is the median x_m . H_m divides the segments to three sets S_ℓ , S_m and S_r . Below the tree corresponding to this partition of segments.

auxiliary data structure T_m which we associate with the root of T . (Structures similar to T_m are built for all internal nodes of T .) T_m is built as follows. H_m cuts naturally each segment in S_m into two pieces. Let C_ℓ be the left pieces of the segments and C_r the right pieces. We build one tree on C_ℓ and one on C_r . We describe the construction of the tree only for C_r since it is symmetric for C_ℓ .

Let x'_m be the median among all x -coordinates of the right endpoints of the segments in C_r . (Note that the x -coordinates of the left endpoints are all equal.) Let H'_m be the plane passing through the point $(x'_m, 0, 0)$ and which is parallel to the yz plane. We store x'_m at the root. The segments of C_r that were not cut by H'_m form the set S'_ℓ . The right pieces of the segments cut by H'_m form the S'_r . The left

pieces (which span between planes H_m and H'_m) form the set S'_m . For S'_ℓ , S'_r we continue the construction recursively (unless empty), much like we built our interval tree T above. See Fig. 3 for an example of such a construction. We use the set of segments S'_m to construct a 2D Voronoi diagram for the point set $H'_m \cap S'_m$. Then this is combined with a standard point location algorithm [4] to give us a data structure T_2 for answering optimally 2D nearest neighbor queries over $H'_m \cap S'_m$. T_2 is associated with the root of the tree for C_r . Structures similar to T_2 are also built for all internal nodes of T_m .

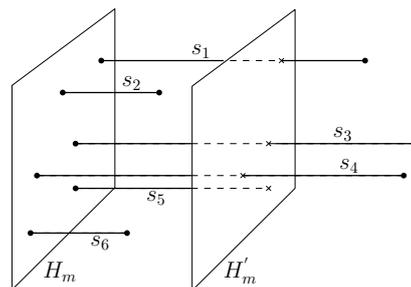


Figure 3: The left endpoint of s_5 is the median x'_m . The parts of segments s_1 , s_2 and s_3 that are on the left of H'_m form the set S'_ℓ . The parts between H_m and H'_m , together with s_5 form S'_m . S'_ℓ is formed by s_2 and s_6 .

We compute a bound on the size of T_m , that is, for the augmented trees on C_ℓ and on C_r . Let $n' = |S_m|$. Because we split always at the median, the height of T_m is $O(\log n')$. This implies that one segment of S_m may be cut at most $O(\log n')$ times and thus the total size of T_m is $O(n' \log n')$. Using standard results, it is easy to see that the total size of all structures T_2 associated with the nodes of T_m is also $O(n' \log n')$.

We compute a bound on the size of the main data structure T . Since the set of segments S_m used at each node of T for the auxiliary data structure T_m are disjoint and using the space bound on T_m it follows easily that the total size of T is $O(n \log n)$. Due to the balanced splits T has height $O(\log n)$.

We describe now how to solve (b) that is, given a query $q = (q_x, q_y, q_z)$, how to find an ε -NN to q in the set $H_q \cap S$. (In fact this first method finds an exact nearest neighbor to q .) We start at the root of T and at each node v we follow the child according to the result of the comparison between q_x and x_m , the value stored at v . At each node v we also visit the auxiliary data structure T_m . We similarly follow the path from the root of T_m to the leaf containing q and at each node we use T_2 to find the nearest neighbor to (q_y, q_z) among $H'_m \cap S'_m$. We report as an answer the nearest point to q over all the points returned from all queries to T_2 .

Correctness follows from the fact that we only ex-

clude from consideration segments or fragments of segments that do not intersect plane H_q and thus are not needed for (b).

Since the depth of tree T is $O(\log n)$ and for each node of T we visit an auxiliary data structure T_m with depth also at most $O(\log n)$ and at each node of T_m the data structure T_2 may have been built for at most n sites, it follows that query time is $O(\log^3 n)$. The total query time is the sum of the times used to solve (a) and (b) and thus we get the following result:

Theorem 1 *Given n parallel segments in 3D we can construct a data structure of $O(n \log n)$ space for finding an ε -NN to any given query point q in $O(\log^3 n + 1/\varepsilon)$ time.*

We will discuss the construction times of the data structures in the full version, however they are all in $O(n \text{ poly}(\log n, \frac{1}{\varepsilon}))$.

2.2 Improving space and query time

We present next a second method that reduces both the space and query time by a $\log n$ factor. The part of the first method that we change is the auxiliary data structure T_m for the sets C_ℓ and C_r . We again describe just the data structure T_r for segments in C_r and an analogous data structure can be built for C_ℓ .

According to (b), our goal for C_r is to find an ε -NN to q in $H_q \cap C_r$. We solve this by reducing our problem to a weight-constrained approximate nearest neighbor searching problem in two dimensions. Specifically each right endpoint (p_x, p_y, p_z) of a segment in C_r is mapped to the point (p_y, p_z) with weight p_x . We denote with P' the 2D weighted point set obtained (see Fig. 4). Let P'_w be the subset of P' containing only the points with weight at least w . Given a query $q = (q_x, q_y, q_z)$, our goal is to find an ε -NN to point $q_2 = (q_y, q_z)$ in P'_w . Note that this suffices to achieve our first goal.

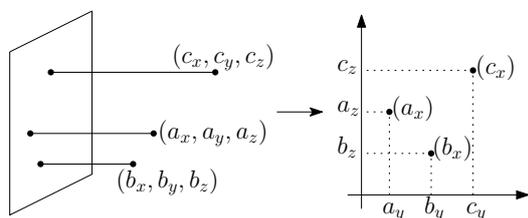


Figure 4: Projection to 2D. The x -coordinate of each point is used as a weight (enclosed in parentheses).

We apply Theorem 4 of the next section (weight-constrained ε -NN queries) on point set P' for $d = 2$, $\gamma = 2$ and $q = q_2$ and easily get this lemma:

Lemma 2 *Given q and C_r we can build a data structure T_r of $O(|C_r| \log(1/\varepsilon))$ size to find an ε -NN to q in $H_q \cap C_r$ in $O(\log |C_r| + 1/\varepsilon^2)$ time.*

Using the above lemma we obtain the following result: (We omit the analysis which is similar to that of the first method.)

Theorem 3 *Given n parallel segments in 3D we can construct a data structure of $O(n \log(1/\varepsilon))$ space for finding an ε -NN to any given query point q in $O(\log^2 n + \log n/\varepsilon^2)$ time.*

3 Weight-constrained ε -NN queries

Given a set of weighted d -dimensional points, we define the *weight-constrained ε -approximate nearest neighbor* problem: given a query q and a weight w , find an ε -NN to q among the points in P that have weight at least w . Here w is a number specified at query time. Note that we allow an approximation error in one parameter (distance) but we require exactness on another (weight). We can also define the symmetric problem where we search for an ε -NN among the points of P with weight at most w .

For the rest of this section we consider only the maximum version of the problem. We state below our result. This result also provides a space-time tradeoff which is controlled by the parameter γ .

Theorem 4 *Let P be a set of n weighted points in \mathbb{R}^d , and let $0 < \varepsilon < 1/2$ and $2 \leq \gamma \leq 1/\varepsilon$ be two real parameters. We can construct a data structure of $O(n\gamma^d \log(1/\varepsilon))$ space that allows us to answer a weight-constrained ε -NN query in time $O(\log(\gamma n) + 1/(\varepsilon\gamma)^d)$.*

For reasons of space we give here only some of the basic ideas and methods that we have used to prove the theorem. We start with some definitions. Let $b(q, r)$ be a ball of radius r centered at point q . Let $b^+(q, r)$ be a ball of radius $(1 + \varepsilon)r$ centered at q . Let $\bar{b}(q, r)$ be the set of points not contained in $b(q, r)$.

Given P , q and r , an *ε -approximate spherical range maximum query* or simply *ε -range maximum query* returns a point in P along with its weight that lies in $b^+(q, r)$ and has weight at least as large as the maximum weight among all points in $b(q, r)$. There are a number of data structures for answering efficiently ε -range maximum queries. Here we will use the data structure in [1]. For $2 \leq \gamma \leq 1/\varepsilon$ it uses $O(n\gamma^d \log(1/\varepsilon))$ and has query time $O(\log(\gamma n) + 1/(\varepsilon\gamma)^{d-1})$. The dependence on ε in query time can be further improved by using known results on approximate idempotent range searching. This implies a similar improvement on the query time of the theorem however we will not discuss these here.

The key idea is to observe that a weight-constrained ε -NN query can be answered with the help of a carefully chosen series of ε -range maximum queries. Assume that there is a method, e.g. [1], that can answer

an ε -range maximum query in $t(n, \varepsilon)$ time. Given a query point q and a weight w , we search for a weight-constrained ε -NN to q . Suppose that we perform an ε -range maximum query with the range $b(q, r)$ for some r of our choice and that it returns a point at distance r' from q with weight w' . Note that $r' \leq (1 + \varepsilon)r$. If $w' \geq w$ this implies we can limit our search for a weight-constrained ε -NN to q among the points in $b(q, r')$. Otherwise if $w' < w$ we are certain that the answer can only be found among the points in $\bar{b}(q, r)$ (since we know from the answer to the ε -NN maximum range query that all points in $b(q, r)$ have weight at most w'). When all points in P have weights less than w there is no weight-constrained ε -NN. To avoid this case we may assume that there is an auxiliary point far enough from P with infinite weight.

We get that after repeating several ε -range maximum queries with the same center q but for different values of the radius r we will have limited our search in an annulus $A = \bar{b}(q, r_1) \cap b(q, r_2)$ for some values r_1, r_2 with $r_1 < r_2$ and we will have found a point p in A satisfying the weight constraint. Observe that if $r_2 \leq (1 + \varepsilon)r_1$ clearly p is a valid answer to the query pair q and w . We show next that with a careful selection of the range radii we can easily obtain an efficient solution for our problem when P has small spread. Let the *spread* Δ of a point set P be the ratio between its diameter and the distance of a closest pair of P . Clearly we can perform a binary search for the right radius using at most $O(\log \Delta)$ ε -range maximum queries and thus find a weight-constrained ε -NN in $O(t(n, \varepsilon) \log \Delta)$ time.

The drawback of the approach described above is that it depends on the spread and it pays at least a $O(\log n)$ factor for each ε -range maximum query even though all of these queries share the same center and may also have similar radii. Interestingly though using together the methods of this approach and the methods presented in [1] for answering ε -approximate range queries and particularly for answering ε -approximate k th nearest neighbor queries we can remove this drawback and arrive at the theorem. Roughly the main adaptation is that wherever an approximate range counting query is needed for approximate k th nearest neighbor searching we use instead the corresponding approximate range maximum query and we guide the search according to the point and the weight returned and not the count. A detailed proof will appear in the full version.

4 Querying about the past

We define a *timestamped* operation on a data structure as an operation which carries a label of the time it occurred. An element is *alive* at some moment t if t is between the time of insertion and deletion of the element. We consider the following problem: given

a sequence of n timestamped insertions and deletions of d -dimensional points build a data structure which given a query point q and a parameter t finds efficiently an ε -NN of q among the points alive at time t . We call this a *t -moment ε -NN query*.

We tackle the problem using methods from Sec. 2 and 3. Let p be a point that was inserted at time t_s and deleted at time t_f (infinite values are allowed). We use time as an extra dimension and map the point p to the segment with endpoints (t_s, p) and (t_f, p) in $d + 1$ dimensions. Note that the points alive at any given moment t correspond to the segments intersecting the hyperplane with equation $x = t$. Hence we can build a similar interval tree as in Sec. 2. To answer queries part (b) is only needed. We extend Lemma 2 in d dimensions and after a simple analysis we get:

Theorem 5 *Let n timestamped insertions and deletions of points in \mathbb{R}^d , and let $0 < \varepsilon < 1/2$ be a real parameter. We can construct a data structure of $O(n \log(1/\varepsilon))$ space that allows us to answer any t -moment ε -NN query in time $O(\log^2 n + \log n/\varepsilon^d)$. Here t is given at query time.*

References

- [1] S. Arya, T. Malamatos, and D. M. Mount. Space-time tradeoffs for approximate spherical range counting. In *Proc. 16th Annu. ACM-SIAM Sympos. Discrete Algorithms*, pages 535–544, 2005.
- [2] S. Arya, T. Malamatos, and D. M. Mount. Space-time tradeoffs for approximate nearest neighbor searching. *J. Assoc. Comput. Mach.*, 57:1–54, 2009.
- [3] S. Arya, D. M. Mount, N. Netanyahu, R. Silverman, and A. Y. Wu. An optimal algorithm for approximate nearest neighbor searching in fixed dimensions. *J. Assoc. Comput. Mach.*, 45:891–923, 1998.
- [4] M. de Berg, O. Cheong, M. van Kreveld, and M. Overmars. *Computational Geometry: Algorithms and Applications*. Springer-Verlag, Berlin, Germany, 3rd edition, 2008.
- [5] P. Indyk. Nearest neighbors in high-dimensional spaces. In *Handbook of Discrete and Computational Geometry*, pages 877–892. CRC Press, 2004.
- [6] P. Indyk and R. Motwani. Approximate nearest neighbors: Towards removing the curse of dimensionality. In *Proc. 30th Annu. ACM Sympos. Theory Comput.*, pages 604–613, 1998.
- [7] V. Koltun and M. Sharir. Polyhedral Voronoi diagrams of polyhedra in three dimensions. In *Proc. 18th Annu. ACM Sympos. Comput. Geom.*, pages 227–236, 2002.
- [8] A. Magen. Dimensionality reductions in ℓ_2 that preserve volumes and distance to affine spaces. *Discrete Comput. Geom.*, 38(1):139–153, 2007.
- [9] Y. Wang. Approximating nearest neighbor among triangles in convex position. *Inf. Process. Lett.*, 108(6):379–385, 2008.