

LINEAR-TIME APPROXIMATE HAUSDORFF DISTANCE

OLIVER CHUBET, PARTH PARIKH, DONALD R. SHEEHY, AND SIDDHARTH SHETH

North Carolina State University

1. INTRODUCTION

The Hausdorff distance is a metric on compact subsets of a metric space. Let (X, \mathbf{d}) be a metric space and let A and B compact subsets of X . The distance from a point $x \in X$ to the set B is $\mathbf{d}(x, B) := \min_{b \in B} \mathbf{d}(x, b)$. The directed Hausdorff distance is $\mathbf{d}_h(A, B) := \max_{a \in A} \mathbf{d}(a, B)$, and the (undirected) Hausdorff distance is $\mathbf{d}_H(A, B) := \max\{\mathbf{d}_h(A, B), \mathbf{d}_h(B, A)\}$. This definition leads directly to a quadratic time algorithm for finite sets. If one only has access to distances (i.e., no ability to hash points), there is a corresponding quadratic-time lower bound.

It is not easy to get an asymptotic improvement on the naive Hausdorff distance algorithm without using an efficient data structure in higher dimensions. Alt et al. [1] give an $O(n \log n)$ time algorithm to compute exact Hausdorff distance in the plane using Voronoi diagrams. Many heuristics are used in practice to speed up the naive algorithm [8, 7, 2]. Another popular technique is to use a geometric tree data structure [9, 5].

We present an algorithm that computes a $(1 + \varepsilon)$ -approximation to the Hausdorff distance between two sets of total size n in $(\frac{1+\varepsilon}{\varepsilon})^{O(d)} n$ time after time after preprocessing the input sets individually into linear-size ball trees. The preprocessing takes $(\frac{1}{\varepsilon})^{O(d)} n \log \Delta$ time per set. Throughout, we assume that distance computations take constant time. Although we explicitly consider the spread of the input in the analysis of our algorithms, there are several cases where we assume that the spread is at most $2^{O(n)}$.

2. BACKGROUND

2.1. Doubling Metrics. Let (X, \mathbf{d}) be a metric space. A **metric ball** in X of radius r centered at c is defined as $\mathbf{ball}(c, r) := \{x \in X \mid \mathbf{d}(x, c) \leq r\}$. The **spread** Δ of $A \subseteq X$ is the ratio of the largest to smallest pairwise distance of points in A .

The set A is **λ -packed** if $\mathbf{d}(a, b) \geq \lambda$ for any distinct $a, b \in A$. A collection of sets \mathcal{C} **covers** A if $A \subseteq \bigcup_{S \in \mathcal{C}} S$. The **doubling constant** of X is the minimum number ρ such that any ball in X can be covered by at most ρ balls of half the radius. The **doubling dimension** is $\dim(X) := \log_2 \rho$. If $\dim(X)$ is bounded then X is a **doubling metric**. The following lemma [4] is true for packed and bounded sets.

Lemma 1 (Standard Packing Lemma). *If X is a metric space with $\dim(X) = d$ and $Z \subset \mathbf{ball}(x, r)$ for some $x \in X$ is λ -packed then $|Z| \leq (\frac{4r}{\lambda})^d$.*

2.2. Greedy Permutations. Let $P = (p_0, \dots, p_{n-1})$ be an ordering of n points. The **i^{th} -prefix** of P is the set P_i containing points p_0, \dots, p_{i-1} . We say P is a **greedy permutation** if $\mathbf{d}(p_i, P_i) = \mathbf{d}_H(p_i, P)$ for all $i > 0$. Let $\alpha > 1$. We say P is an α -approximate greedy permutation if $\mathbf{d}_H(p_i, P) \leq \alpha \mathbf{d}(p_i, P_i)$ for all $i > 0$. If $\mathbf{d}(p_i, q) \leq \alpha \mathbf{d}(p_i, P_i)$, then q is an **approximate nearest predecessor**. We denote the distance to the

(approximate) nearest predecessor of p_i by ε_{p_i} . Then there is a $\frac{1}{\alpha}\varepsilon_{p_i}$ -packing of the points in the prefix P_i . The greedy permutation of a set can be computed in $O(n \log \Delta)$ time in low dimensions.

2.3. Greedy Trees. A balltree [6] is defined by recursively partitioning a metric space and representing the parts in a binary tree. Each node of the tree is a metric ball that covers the points in its subtree.

A **greedy tree** G is a balltree that uses the greedy permutation P to guide the partition. For each node $\mathbf{ball}(p, r)$, the center p is a point from P . The radius r is the maximum distance to a point of P in the ball. Any node $\mathbf{ball}(p, r)$ with $r > 0$ has two children, $\mathbf{ball}(p, r_L)$ and $\mathbf{ball}(c, r_R)$, where p is the (approximate) nearest predecessor of c .¹ The nodes with radius 0 are leaves; there is a node $\mathbf{ball}(p, 0)$ for each element of P . The greedy tree has $2n - 1$ nodes, where $|P| = n$, and the radii are non-increasing from parent to child. Given a greedy permutation and the approximate nearest predecessors, the corresponding greedy tree can be computed in $O(n \log \Delta)$ time ($O(n)$ time to build the tree and $O(n \log \Delta)$ time to compute radii).

3. APPROXIMATE HAUSDORFF DISTANCE

In this section, we describe an algorithm that computes $d_h(A, B)$, given sets A and B . As in prior work [3], we assume that the input is preprocessed into a tree structure; in our case, we use greedy trees. We assume that nodes of the greedy trees are listed in sorted order by radii. The algorithm proceeds by iterating over the nodes. Each node gets added as a vertex in the **neighbor graph** which maintains the following invariants:

- **Neighbor Invariant:** If $d(a, B) = d(a, b)$ then there is an edge between the nodes storing a and b .
- **Partition of Input:** Any point in A or B will be stored by some node in the neighbor graph.
- **Lower Bounds:** A node p in the neighbor graph stores a lower bound $l(p)$ to the distance $d(p, B)$.

This is called a **local lower bound**. The greatest of these is the **global lower bound**, denoted L .

The algorithm then proceeds, updating the neighbor graph and pruning edges that are too long to impact the neighbor invariant. It stops when it reaches a node whose radius is sufficiently small compared to the global lower bound. At that point the lower bound is a good approximation, and we return it. We update

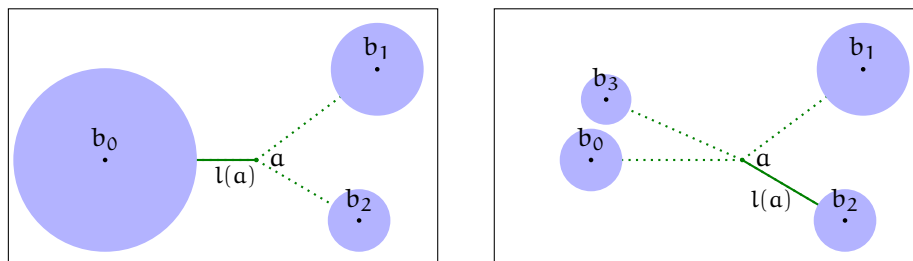


FIGURE 1. This figure depicts an update of $l(a)$ after replacing a node at b_0 with its children.

the local lower bound for a point p as shown in Figure 1 by computing,

$$l(p) = \min_{q \in N(p)} \{d(p, q) - r_q\}.$$

We update L each time we update a local lower bound by comparing the two values (see Figure 2). We prune nodes that are too far to contain any nearest neighbors (see Section 3.1) and stop when L is a $(1 + \varepsilon)$ -approximation of the exact distance (see Section 3.2). These improvements bound degrees in the neighbor graph by a constant and so each node can be processed in constant time.

¹The (approximate) nearest predecessors need not be unique, however for the sake of construction we assume we have chosen one.

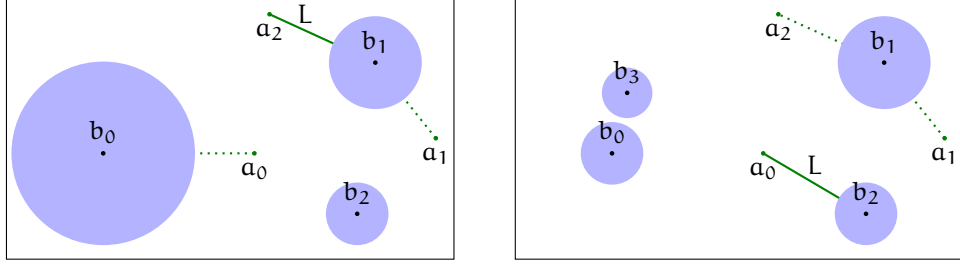


FIGURE 2. This figure depicts an update of L after replacing the node at b_0 with its children.

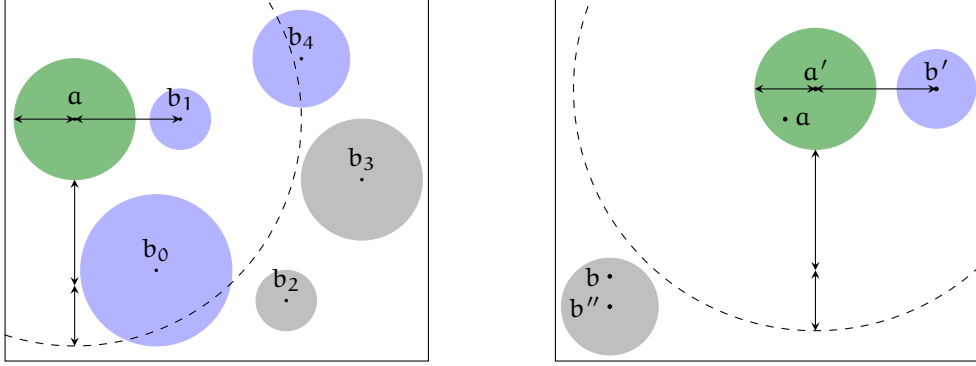


FIGURE 3. This figure shows the role of the pruning condition. On the left, the balls b_2 and b_3 are too far away to contain the nearest neighbor of a , so edges (a, b_2) and (a, b_3) will be pruned from the neighbor graph. The pruning condition respects the neighbor invariant and does not prune edge (a, b_4) . In the right image let $a' \in A^r$. For any point $a \in B(a', r_{a'})$, by the triangle inequality, $d(a, B) \leq r_{a'} + d(a', b')$. If b'' is a pruned center, then no point $b \in B(b'', r_{b''})$ can be the nearest neighbor of a , because $d(a, b') \leq d(a, b)$ for any such b .

Algorithm: HAUSDORFF

Input: Greedy trees G_A, G_B and approximation factor ε

Output: $(1 + \varepsilon)$ -approximation of $d_h(A, B)$

- (1) Initialize the neighbor graph N with a node for each root and an edge between them.
- (2) Iterate over nodes $\mathbf{ball}(p, r)$ of $G_A \cup G_B$ in non-increasing order of radii while $r > (\frac{\varepsilon}{2})L$.
 - (a) Let the children of $\mathbf{ball}(p, r)$ be called $\mathbf{ball}(p, r_L)$ and $\mathbf{ball}(c, r_R)$.
 - (b) Add c to N , and add edges so that $N(c) = N(p)$.
 - (c) Update the radii, $r_p = r_L$ and $r_c = r_R$.
 - (d) If $p \in A$, prune $N(c)$ and $N(p)$, and update $l(c)$ and $l(p)$.
 - (e) If $p \in B$, for each $a \in N(p)$, prune $N(a)$ and update $l(a)$.
- (3) Return L .

To compute exact $d_h(A, B)$, let $\varepsilon = 0$. However our running time guarantees do not hold in that case.

3.1. Pruning Condition. Consider an iteration of the main loop where the node $\mathbf{ball}(p, r)$ is being processed. Let A^r and B^r be the nodes in N from G_A and G_B respectively. We prune an edge $(a, b) \in N$ if $d(a, b') + r_a < d(a, b) - r_a - r_b$ for some $b' \in B^r$, as shown in Figure 3. The following lemma proves that the pruning condition does not trim any edge between nodes that contain nearest neighbors. A proof of correctness for the pruning condition is presented in Appendix B.

3.2. Stopping Condition. The running time of an iteration depends on the degree of a vertex in the neighbor graph. The key to achieving this bound is to stop the algorithm early so that the number of neighbors can be bounded by a packing argument. Let r be the radius of the current node, and L the global lower bound. Then, $L \leq \mathbf{d}_h(A, B) \leq L + 2r$. Moreover, if $r \leq (\frac{\epsilon}{2})L$, then L is a $(1 + \epsilon)$ -approximation to $\mathbf{d}_h(A, B)$ (see Appendix B).

3.3. Analysis.

Theorem 2. *Given two approximate greedy trees for sets A and B of total cardinality n , HAUSDORFF computes a $(1 + \epsilon)$ -approximation of $\mathbf{d}_h(A, B)$ in $(\frac{1+\epsilon}{\epsilon})^{O(d)} n$ time.*

Proof. In order to bound the degrees of the neighbor graph N , we first establish that the points associated with the neighbors of a vertex in N are packed. By construction, any node $p \in N$ is the center of the left- or right-child of a greedy tree node with radius at least r . Then by Lemma 3, N is $\frac{(\alpha-1)r}{(2\alpha-1)\alpha}$ -packed. Thus,

$$|N(a)| \leq \left(\frac{2\alpha(2\alpha-1)(L+2r)}{(\alpha-1)r} \right)^d \leq \left(\frac{8\alpha^2(1+\epsilon)}{(\alpha-1)\epsilon} \right)^d,$$

for all $r \geq (\frac{\epsilon}{2})L$, by Lemma 1. Therefore, the number of edges incident to any given node in N is $(\frac{1+\epsilon}{\epsilon})^{O(d)}$. So we spend $(\frac{1+\epsilon}{\epsilon})^{O(d)}$ time for each iteration of the algorithm. This gives a running time of $(\frac{1+\epsilon}{\epsilon})^{O(d)} n$. \square

4. CONCLUSION

We have given a new algorithm for computing the Hausdorff distance and its relatives. After pre-processing the point sets, individual distance computations take only linear time.

REFERENCES

- [1] Helmut Alt, Bernd Behrends, and Johannes Blömer. Approximate matching of polygonal shapes. *Annals of Mathematics and Artificial Intelligence*, 13(3-4):251–265, September 1995.
- [2] Yilin Chen, Fazhi He, Yiqi Wu, and Neng Hou. A local start search algorithm to compute exact Hausdorff Distance for arbitrary point sets. *Pattern Recognition*, 67:139–148, July 2017.
- [3] Daniel P. Huttenlocher, Klara Kedem, and Jon M. Kleinberg. On dynamic voronoi diagrams and the minimum hausdorff distance for point sets under euclidean motion in the plane. In *Proceedings of the Eighth Annual Symposium on Computational Geometry*, SCG '92, pages 110–119, New York, NY, USA, 1992. Association for Computing Machinery.
- [4] Robert Krauthgamer and James R. Lee. Navigating nets: Simple algorithms for proximity search. In *Proceedings of the Fifteenth Annual ACM-SIAM Symposium on Discrete Algorithms*, SODA '04, pages 798–807, USA, 2004. Society for Industrial and Applied Mathematics.
- [5] Sarana Nutanong, Edwin H. Jacox, and Hanan Samet. An incremental Hausdorff distance calculation algorithm. *Proceedings of the VLDB Endowment*, 4(8):506–517, May 2011.
- [6] Stephen M. Omohundro. Five balltree construction algorithms. Technical Report 562, ICSI Berkeley, 1989.
- [7] Jegoon Ryu and Sei-ichiro Kamata. An efficient computational algorithm for Hausdorff distance based on points-ruling-out and systematic random sampling. *Pattern Recognition*, 114:107857, June 2021.
- [8] Abdel Aziz Taha and Allan Hanbury. An Efficient Algorithm for Calculating the Exact Hausdorff Distance. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(11):2153–2163, November 2015.
- [9] Dejun Zhang, Fazhi He, Soonhung Han, Lu Zou, Yiqi Wu, and Yilin Chen. An efficient approach to directly compute the exact Hausdorff distance for 3D point sets. *Integrated Computer-Aided Engineering*, 24(3):261–277, July 2017.

APPENDIX A. PACKING BY GREEDY TREES

Here we present some properties of the greedy tree that allow us to bound the degree of a vertex in the neighbor graph at any stage of HAUSDORFF.

Lemma 3. *Let G_A be an α -approximate greedy tree and let $X \subseteq G_A$ be the set of all nodes in G_A with radius at least r . Then the centers of all children of the nodes of X are $\frac{(\alpha-1)r}{(2\alpha-1)\alpha}$ -packed.*

Proof. We can establish a packing by finding a lower bound of the insertion distances. We note that each center of a left-child (left-center) of X is a predecessor to some center of a right-child (right-center) of X . Therefore every left-center has an insertion distance greater than some right-center of X . Thus, it is sufficient to only consider the right-centers of X .

Let $B(p, r_p) \in X$ have a right-child $B(c, r_c)$. Let C be the set of c' such that p is the predecessor of c' . To bound ε_c , first consider r_p . We know $r_p \leq \max_{c' \in C} \{\varepsilon_{c'} + r_{c'}\}$. From the radius bound it follows $r_p \leq \max_{c' \in C} \{\varepsilon_{c'} + (\frac{\alpha}{\alpha-1})r_{c'}\}$. The greedy tree construction guarantees that $\max_{c' \in C} \{\varepsilon_{c'}\} = \varepsilon_c$, so $r_p \leq (\frac{2\alpha-1}{\alpha-1})\varepsilon_c$. This gives a lower bound on the insertion distance,

$$\varepsilon_c \geq \left(\frac{\alpha-1}{2\alpha-1}\right) r_p \geq \left(\frac{\alpha-1}{2\alpha-1}\right) r.$$

This bound holds for any right-center of X , therefore the centers of the children of X are $\frac{(\alpha-1)r}{(2\alpha-1)\alpha}$ -packed. \square

Lemma 4. *For any greedy tree node $B(q, r_q)$ we have $r_q \leq (\frac{\alpha}{\alpha-1})\varepsilon_q$.*

Proof. In an α -approximate greedy tree, $\varepsilon_c \leq \frac{1}{\alpha}\varepsilon_q$ where q is the nearest predecessor of c . It follows,

$$r_q \leq \sum_{i=1}^{\infty} \left(\frac{1}{\alpha}\right)^i \varepsilon_q \leq \left(\frac{\alpha}{\alpha-1}\right)\varepsilon_q.$$

\square

APPENDIX B. PRUNING AND STOPPING IN HAUSDORFF

The following lemma proves the correctness of the pruning condition of HAUSDORFF.

Lemma 5. *Let N be the neighbor graph at radius r . For any $a \in A$, if $\mathbf{d}(a, B) = \mathbf{d}(a, b)$ then there exists an edge (a', b') in N such that $a \in \mathbf{ball}(a', r_{a'})$ and $b \in \mathbf{ball}(b', r_{b'})$.*

Proof. Given $a \in A$, let $b \in B$ such that $\mathbf{d}(a, b) = \mathbf{d}(a, B)$. The neighbor graph maintains a partition of A and B , so there exists $a' \in A^r$ and $b' \in B^r$ such that $a \in \mathbf{ball}(a', r_{a'})$ and $b \in \mathbf{ball}(b', r_{b'})$. Then, $\mathbf{d}(a, b) \leq \mathbf{d}(a', b') - r_{a'} - r_{b'}$ by the triangle inequality. Suppose in the current iteration we pruned the edge (a', b') in N . Then (a', b') must satisfy the pruning condition, $\mathbf{d}(a', b'') + r_{a'} < \mathbf{d}(a', b') - r_{a'} - r_{b'}$, for some $b'' \in B^r$. It follows that $\mathbf{d}(a', b'') + r_{a'} < \mathbf{d}(a, b)$. Yet, $\mathbf{d}(a, b) \leq r_{a'} + \mathbf{d}(a', b'')$, so this is a contradiction. Therefore, we cannot prune an edge between nodes storing nearest neighbors. \square

The following lemma proves that L is a $(1 + \varepsilon)$ -approximation of $\mathbf{d}_h(A, B)$ when r is sufficiently small in HAUSDORFF.

Lemma 6. *Let r be the radius of the current node, and L the global lower bound. Then, $L \leq \mathbf{d}_h(A, B) \leq L + 2r$. Moreover, if $r \leq (\frac{\varepsilon}{2})L$, then L is a $(1 + \varepsilon)$ -approximation of $\mathbf{d}_h(A, B)$.*

Proof. We first show that the distance from A^r to B is at most $L + r$. We know that $\mathbf{d}_h(A^r, B) \leq \mathbf{d}_h(A^r, B^r)$ because $\mathbf{d}(a, B) \leq \mathbf{d}(a, B^r)$ for any $a \in A^r$. We also know that $r_p \leq r$ for any node $p \in N$. So,

$$\mathbf{d}(a, B^r) \leq \min_{b \in B^r} \{\mathbf{d}(a, b) - r_b + r\} = l(a) + r.$$

It follows that $\mathbf{d}(A^r, B) \leq \max_{a \in A^r} \{l(a) + r\} = L + r$. Furthermore, $\mathbf{d}(a, B) \leq \mathbf{d}(a, A^r) + \mathbf{d}_h(A^r, B)$, and we know that $\mathbf{d}(a, A^r) \leq r$. So, $\mathbf{d}(a, B) \leq L + 2r$. Therefore $L \leq \mathbf{d}_h(A, B) \leq L + 2r$. \square