



The VC Dimension of Metric Balls under Fréchet and Hausdorff Distances

Anne Driemel¹ · André Nusser² · Jeff M. Phillips³ · Ioannis Psarros^{1,4} 

Received: 21 November 2019 / Revised: 19 October 2020 / Accepted: 8 February 2021

© The Author(s) 2021

Abstract

The Vapnik–Chervonenkis dimension provides a notion of complexity for systems of sets. If the VC dimension is small, then knowing this can drastically simplify fundamental computational tasks such as classification, range counting, and density estimation through the use of sampling bounds. We analyze set systems where the ground set X is a set of polygonal curves in \mathbb{R}^d and the sets \mathcal{R} are metric balls defined by curve similarity metrics, such as the Fréchet distance and the Hausdorff distance, as well as their discrete counterparts. We derive upper and lower bounds on the VC dimension that imply useful sampling bounds in the setting that the number of curves is large, but the complexity of the individual curves is small. Our upper and lower bounds are either near-quadratic or near-linear in the complexity of the curves that define the ranges and they are logarithmic in the complexity of the curves that define the ground set.

Keywords VC dimension · Polygonal curves · Fréchet distance · Hausdorff distance

Mathematics Subject Classification 68P01 · 68U05 · 68R01

Editor in Charge: Kenneth Clarkson

We thank Peyman Afshani for useful discussions on the topic of this paper. We also thank the organizers of the 2016 NII Shonan Meeting “Theory and Applications of Geometric Optimization” where this research was initiated. Anne Driemel thanks the Hausdorff Center for Mathematics for their generous support and the Netherlands Organization for Scientific Research (NWO) for support under Veni Grant 10019853. Jeff Phillips thanks for his support from the National Science Foundation (NSF) through Grants CCF-1350888, ACI-1443046, CNS-1514520, CNS-1564287, and IIS-1816149. Part of the work was completed while visiting the Simons Institute for Theory of Computing. Ioannis Psarros thanks the State Scholarships Foundation: this research is co-financed by Greece and the European Union (European Social Fund—ESF) through the Operational Programme « Human Resources Development, Education and Lifelong Learning » in the context of the project “Strengthening Human Resources Research Potential via Doctorate Research” (MIS-5000432), implemented by the State Scholarships Foundation (IKY).

Extended author information available on the last page of the article

1 Introduction

A *range space* (X, \mathcal{R}) (also called *set system*) is defined by a ground set X and a set of ranges \mathcal{R} , where each $r \in \mathcal{R}$ is a subset of X . A data structure for range searching answers queries for the subset of the input data that lies inside the query range. In range counting, we are interested only in the size of this subset. In our setting, a range is a metric ball defined by a curve and a radius. The ball contains all curves that lie within this radius from the center under a specific distance function (e.g., Fréchet or Hausdorff distance).

A crucial descriptor of any range space is its VC dimension [41,43,46] and related shattering dimension, which we define formally below. These notions quantify how complex a range space is, and have played fundamental roles in machine learning [7, 45], data structures [17], and geometry [14,30]. For instance, specific bounds on these complexity parameters are critical for tasks as diverse as neural networks [7,36], art-gallery problems [26,37,44], and kernel density estimation [35].

The last five years have seen a surge of interest in data structures for trajectory processing under the Fréchet distance, manifested in a series of publications [2,8–11,15,21,22,24,29,47]. This was partially motivated by the increasing availability and quality of trajectory data from mobile phones, GPS sensors, RFID technology, and video analysis [28,38,48]. Initial results in this line of research, such as the approximate range counting data structure by de Berg et al. [10], use classical data structuring techniques. Afshani and Driemel extended their results and in addition showed lower bounds on the space-query-time trade-off in this setting [2]. In particular, they showed a lower bound which is exponential in the complexity of the curves for exact range searching. In 2017, ACM SIGSPATIAL, the premier conference for geographic information science, devoted their software challenge (GIS CUP) to the problem of range searching under the Fréchet distance [47]. Spurring further developments, the most recent results explore the use of heuristics [13] and randomization [16].

The Fréchet distance, named after Maurice Fréchet [25], is a popular distance measure for curves. Intuitively, it can be defined using the metaphor of a person walking a dog, where the person follows one curve and the dog follows the other curve, and throughout their traversal they are connected by a leash of fixed length. Both can vary their speed but they are not allowed to move backwards. The Fréchet distance corresponds to the length of the shortest dog leash that permits a traversal in this fashion. The Fréchet distance is very similar to the Hausdorff distance for sets [32], which is defined as the minimal maximum distance of a pair of points, one from each set, under all possible mappings between the two sets. The difference between the two distance measures is that the Fréchet distance requires the mapping to adhere to the ordering of the points along the curve. Both distance measures allow flexible associations between parts of the input elements which sets them apart from classical L_p distances and makes them so suitable for trajectory data under varying speeds. One standard tool for computing the Fréchet distance of two curves is the free-space diagram which was introduced by Alt and Godau [6]. In the free-space diagram, we consider the polygonal curves as continuous curves $[0, 1] \rightarrow \mathbb{R}^d$. The free-space for a given distance threshold ρ is a subset of the parametric space $[0, 1] \times [0, 1]$ that consists of all point pairs on the two curves at distance at most ρ . The vertices of the

curves partition $[0, 1] \times [0, 1]$ into rectangular cells, such that each cell corresponds to the parametric space of two edges, one from each curve. One can decide if the distance between two polygonal curves is at most ρ by checking whether there is a path which is monotone in both coordinates that starts at $(0, 0)$, ends at $(1, 1)$, and stays inside the free-space.

Our contribution in this paper is a comprehensive analysis of the Vapnik–Chervonenkis dimension of the corresponding range spaces. The resulting VC dimension bounds, while being interesting in their own right, have a plethora of applications through the implied sampling bounds. We detail a range of implications of our bounds in Sect. 10.

2 Definitions

In this section, we formally define the distances between curves as well as VC dimension and range spaces, so we can state our main results. This basic set up will be enough to prove our results for the discrete variants of the distance measures we consider. The basic proofs in the discrete setting also serve as a template for the proofs in the main part of the paper. Starting in Sect. 6 we provide more advanced geometric definitions and properties about the VC dimension which we then use in our proofs on the continuous variants of the distance measures we consider.

2.1 Distance Measures

The Fréchet distance was first defined by Maurice Fréchet in his doctoral thesis of 1906 [25]. The Hausdorff distance was first defined by Felix Hausdorff in his book “Grundzüge der Mengenlehre” of 1914 [32]. Here, we follow the definitions given by Alt and Godau [6] for the continuous variants of the Fréchet distance, we follow Eiter and Mannila [23] for the discrete variant, and we use the original definitions for the Hausdorff distance. We denote by $\|\cdot\|$ the Euclidean norm $\|\cdot\|_2$.

Definition 2.1 (*directed Hausdorff distance*) Let X, Y be two subsets of some metric space (M, d) . The directed Hausdorff distance from X to Y is

$$d_{\vec{H}}(X, Y) = \sup_{u \in X} \inf_{v \in Y} d(u, v).$$

Definition 2.2 (*Hausdorff distance*) Let X, Y be two subsets of some metric space (M, d) . The Hausdorff distance between X and Y is

$$d_H(X, Y) = \max \{d_{\vec{H}}(X, Y), d_{\vec{H}}(Y, X)\}.$$

Definition 2.3 Given polygonal curves V and U with vertices v_1, \dots, v_{m_1} and u_1, \dots, u_{m_2} respectively, a traversal $T = (i_1, j_1), \dots, (i_t, j_t)$ is a sequence of pairs of indices referring to a pairing of vertices from the two curves such that:

1. $i_1, j_1 = 1, i_t = m_1, j_t = m_2$;

2. $\forall (i_k, j_k) \in T: i_{k+1} - i_k \in \{0, 1\}$ and $j_{k+1} - j_k \in \{0, 1\}$;
3. $\forall (i_k, j_k) \in T: (i_{k+1} - i_k) + (j_{k+1} - j_k) \geq 1$.

Definition 2.4 (*discrete Fréchet distance*) Given polygonal curves V and U with vertices v_1, \dots, v_{m_1} and u_1, \dots, u_{m_2} respectively, we define the discrete Fréchet distance between V and U as the following function:

$$d_{dF}(V, U) = \min_{T \in \mathcal{T}} \max_{(i_k, j_k) \in T} \|v_{i_k} - u_{j_k}\|,$$

where \mathcal{T} denotes the set of all possible traversals for V and U .

Any polygonal curve V with vertices v_1, \dots, v_{m_1} and edges $\overline{v_1 v_2}, \dots, \overline{v_{m_1-1} v_{m_1}}$ has a uniform parametrization that allows us to view it as a parametrized curve $v: [0, 1] \rightarrow \mathbb{R}^2$. In the remainder, we use the term curve to refer to polygonal curves if not mentioned otherwise.

Definition 2.5 (*Fréchet distance*) Given two curves $u, v: [0, 1] \rightarrow \mathbb{R}^2$, their Fréchet distance is defined as follows:

$$d_F(u, v) = \min_{\substack{f: [0,1] \rightarrow [0,1] \\ g: [0,1] \rightarrow [0,1]}} \max_{\alpha \in [0,1]} \|v(f(\alpha)) - u(g(\alpha))\|,$$

where f and g range over all continuous, non-decreasing functions with $f(0) = g(0) = 0$ and $f(1) = g(1) = 1$.

Definition 2.6 (*weak Fréchet distance*) Given parametrized curves $u, v: [0, 1] \rightarrow \mathbb{R}^2$, their weak Fréchet distance is defined as follows:

$$d_{wF}(u, v) = \min_{\substack{f: [0,1] \rightarrow [0,1] \\ g: [0,1] \rightarrow [0,1]}} \max_{\alpha \in [0,1]} \|v(f(\alpha)) - u(g(\alpha))\|,$$

where f and g range over all continuous functions with $f(0) = g(0) = 0$ and $f(1) = g(1) = 1$.

2.2 Range Spaces

Each range space can be defined as a pair of sets (X, \mathcal{R}) , where X is the *ground set* and $\mathcal{R} \subseteq 2^X$ is the *range set*. Let (X, \mathcal{R}) be a range space. For $Y \subseteq X$, we denote

$$\mathcal{R}|_Y = \{R \cap Y \mid R \in \mathcal{R}\}.$$

If $\mathcal{R}|_Y$ contains all subsets of Y , then Y is *shattered* by \mathcal{R} .

Definition 2.7 (*VC dimension*) The Vapnik–Chervonenkis dimension [41,43,46] (VC dimension) of (X, \mathcal{R}) is the maximum cardinality of a shattered subset of X .

Definition 2.8 (*shattering dimension*) The shattering dimension of (X, \mathcal{R}) is the smallest δ such that, for all m ,

$$\max_{\substack{B \subset X \\ |B|=m}} |\mathcal{R}|_B = O(m^\delta).$$

It is well known that for a range space (X, \mathcal{R}) with VC dimension ν and shattering dimension δ that $\nu \leq O(\delta \log \delta)$ and $\delta = O(\nu)$. So bounding the shattering dimension and bounding the VC dimension are asymptotically equivalent within a log factor. For a proof of this and other basic facts on range spaces we refer the reader to the textbook of Har-Peled [30].

Definition 2.9 (*dual range space*) Given a range space (X, \mathcal{R}) , for any $p \in X$ we define

$$\mathcal{R}_p = \{R \mid R \in \mathcal{R}, p \in R\}.$$

The dual range space of (X, \mathcal{R}) is the range space $(\mathcal{R}, \{\mathcal{R}_p \mid p \in X\})$.

It is a well-known fact that if a range space has VC dimension ν , then the dual range space has VC dimension $\leq 2^{\nu+1}$ (see e.g. [30]).

There are many techniques for bounding the VC dimension of geometric range spaces. For instance, when the ground set is \mathbb{R}^d and the ranges are defined by inclusion in halfspaces, then the range space and its dual range space are isomorphic and both have VC dimension and shattering dimension d . When the ranges are defined by inclusion in balls, then the VC dimension and shattering dimension is $d + 1$, and the dual range spaces have bounds of d [30]. It is also, for instance, known [12] that the composition ranges formed as the k -fold union or intersection of ranges from a range space with bounded VC dimension ν induces a range space with VC dimension $O(\nu k \log k)$, and it was recently shown by Csikós et al. that this is tight even for some simple range spaces such as those defined by halfspaces [18,19]. More such results are deferred to Sect. 6.

2.3 Range Spaces Induced by Distance Measures

Let (M, d) be a pseudometric space. We define the *ball* of radius r and center p , under the distance measure d , as the following set:

$$b_d(p, r) = \{x \in M \mid d(x, p) \leq r\},$$

where $p \in M$. The doubling dimension of a metric space (M, d) , denoted as $\text{ddim}(M, d)$, is the smallest integer t such that any ball can be covered by at most 2^t balls of half the radius.

In this paper, we study the VC dimension of variants of range spaces (X, \mathcal{R}) induced by pseudometric spaces¹ (M, d) by setting $X = M$ and

$$\mathcal{R} = \{b_d(p, r) \mid r \in \mathbb{R}, r > 0, p \in M\}.$$

It is a reasonable question to ask whether the doubling dimension of a metric space influences the VC dimension of the induced range space. In general, a bounded doubling dimension does not imply a bounded VC dimension of the induced range space and vice versa. Recently, Huang et al. [34] showed that if we allow a small $(1 + \varepsilon)$ -distortion of the distance function d , the shattering dimension can be upper bounded by $O(\varepsilon^{-O(\text{ddim}(M, d))})$. It is conceivable that the doubling dimension of the metric space of the discrete Fréchet distance and Hausdorff distance is bounded, as long as the underlying metric has bounded doubling dimension. However, for the continuous Fréchet distance, the doubling dimension is known to be unbounded [20]. Moreover, we will see that much better bounds can be obtained by a careful study of the specific distance measure.

Specifically, we study an *unbalanced* version of the above range space, in the sense that we distinguish between the complexity of objects of the ground set and the complexity of objects defining the ranges. In our case, the ground set consists of polygonal curves of complexity m , and the ranges are defined by polygonal curves of complexity k . To this end, we define, for any integers d and m , $\mathbb{X}_m^d := (\mathbb{R}^d)^m$ and we treat the elements of this set as ordered sets of points in \mathbb{R}^d of size m . Formally, we study range spaces with ground set \mathbb{X}_m^d and a range set of the form

$$\mathcal{R}_{d,k} = \{b_d(p, r) \cap \mathbb{X}_m^d \mid r \in \mathbb{R}, r > 0, p \in \mathbb{X}_k^d\}$$

under different variants of the Fréchet and Hausdorff distances. We emphasize that the range space consists of ranges of all radii.

3 Our Results

Table 1 shows an overview of our bounds. For metric balls defined on point sets (resp. point sequences) in \mathbb{R}^d we show that the VC dimension is at most near-linear in dk , the complexity of the ball centers that define the ranges, and at most logarithmic in dm , the complexity of point sets of the ground set. Our lower bounds show that these bounds are almost tight in all parameters k , d , and m . For the Hausdorff distance, where the ground set X consists of continuous polygonal curves in \mathbb{R}^d , we show an upper bound that is quadratic in k , quadratic in d , and logarithmic in m . The same bound holds for the Fréchet distance, where the ground set consists of sets of line segments in \mathbb{R}^d . We obtain slightly better bounds in k for the weak Fréchet distance. Our lower bounds extend to the continuous case, but are only tight in the dependence on m – the complexity of the ground set.

¹ While we may use the term *metric* or *pseudometric* to define the range, our methods do not assume any metric properties of the inducing distance measure.

Table 1 Our bounds on the VC dimension of range spaces of the form $(\mathbb{X}_m^d, \mathcal{R}_{d,k})$, for d being the distance measures in the table

Disc.	Hausdorff	$O(dk \log dkm)$ (Theorems 5.1, 5.2)	$\Omega(\max(dk \log k, \log dm))$
	Fréchet		$(d \geq 4, \text{Theorem 9.7})$
Cont.	Hausdorff	$O(d^2 k^2 \log dkm)$ (Theorem 7.7)	$\Omega(\max(k, \log m))$
	weak Fréchet	$O(d^2 k \log dkm)$ (Theorem 8.3)	$(d \geq 2, \text{Theorem 9.4})$
	Fréchet	$O(d^2 k^2 \log dkm)$ (Theorem 8.5)	

In the first column we distinguish between \mathbb{X}_m^d consisting of *discrete* point sequences vs. \mathbb{X}_m^d consisting of *continuous* polygonal curves. The lower bounds hold for all distance measures in this table

While the VC dimension bounds for the discrete Hausdorff and Fréchet metric balls may seem like an easy implication of composition theorems for the VC dimension [12, 18], we still find three things about these results remarkable:

1. First consider the valid alignment paths in the free-space diagram: those are all sequences of cells which are monotonic in both coordinates, their first cell contains $(0, 0)$, and their last cell contains $(1, 1)$. For Fréchet variants, there are $\Theta(2^k 2^m)$ valid alignment paths in the free-space diagram. And one may expect that these may materialize in the size of the composition theorem. Yet by a simple analysis of the shattering dimension, we show that they do not.
2. Second, the VC dimension only has logarithmic dependence on the size m of the curves in the ground set, rather than a polynomial dependence one would hope to obtain by simple application of composition theorems. This difference has important implications in analyzing real data sets where we can query with simple curves (small k), but may not have a small bound on the size of the curves in the data set (large m).
3. Third, for the continuous variants, the range spaces can indeed be decomposed into problems with ground sets defined on line segments. However, we do not know of a general d -dimensional bound on the VC dimension of range space with a ground set of segments, and ranges defined by segments within a radius r of another segment. We are able to circumvent this challenge with a technique to bound the VC dimension using a simple model of computation, and careful predicate design.

4 Our Approach

Our methods use the fact that both the Fréchet distance and the Hausdorff distance are determined by one of a discrete set of events, where each event involves a constant number of simple geometric objects. For example, it is well known that the Hausdorff distance between two discrete sets of points is equal to the distance between two points from the two sets. The corresponding event happens as we consider a value $\delta > 0$ increasing from 0 and we record which points of one set are contained in which balls of radius δ centered at points from the other set. The same phenomenon is true for the discrete Fréchet distance between two point sequences. In particular, the so-called free-space matrix (the discrete version of the free-space diagram) which can be used

to decide whether the discrete Fréchet distance is smaller than a given value δ encodes exactly the information about which pairs of points have distance at most δ . The basic phenomenon remains true for the continuous versions of the two distance measures if we extend the set of simple geometric objects to include line segments and if we also consider triple intersections. Each type of event can be translated into a range space of which we can analyze the VC dimension. Together, the product of the range spaces encodes the information, which curves lie inside which metric balls, in the form of a set system. This representation allows us to prove bounds on the VC dimension of metric balls under these distance measures.

5 Basic Idea: Discrete Fréchet and Hausdorff

In this section we prove our upper bounds in the discrete setting. Let $\mathbb{X}_m^d = (\mathbb{R}^d)^m$; we treat the elements of this set as ordered sets of points in \mathbb{R}^d of size m . The range spaces that we consider in this section are defined over the ground set \mathbb{X}_m^d and the range set of balls under either the Hausdorff or the discrete Fréchet distance. The proofs in the subsequent sections all follow the basic idea of the proof in the discrete setting.

Theorem 5.1 *Let $(\mathbb{X}_m^d, \mathcal{R}_{H,k})$ be the range space with $\mathcal{R}_{H,k}$ being the set of all balls under the Hausdorff distance centered at point sets in \mathbb{X}_k^d . The VC dimension is $O(dk \log dkm)$. The shattering dimension is $O(dk \log m)$.*

Proof Let $\{S_1, \dots, S_t\} \subseteq \mathbb{X}_m^d$ and $S = \bigcup_i S_i$; we define S so that it ignores the ordering within each S_i and is a single set of size tm . Any intersection of a Hausdorff ball with $\{S_1, \dots, S_t\}$ is uniquely defined by a set $\{B_1 \cap S, \dots, B_k \cap S\}$, where B_1, \dots, B_k are balls in \mathbb{R}^d . To see that, notice that the discrete Hausdorff distance between two sets of points is uniquely defined by the distances between points of the two sets.

Consider the range space $(\mathbb{R}^d, \mathcal{B})$, where \mathcal{B} is the set of balls in \mathbb{R}^d . It is well known that the VC dimension is $d + 1$. Hence,

$$\max_{S \subseteq \mathbb{R}^d, |S|=tm} |\mathcal{B}|_S = O((tm)^{d+1}).$$

This implies that

$$|\{\{B_1 \cap S, \dots, B_k \cap S\} \mid B_1, \dots, B_k \text{ are balls in } \mathbb{R}^d\}| \leq O((tm)^{(d+1)k}),$$

and hence,²

$$2^t \leq O((tm)^{(d+1)k}) \implies t = O(dk \log dkm).$$

We can similarly bound the shattering dimension δ ,

$$t^\delta \leq O((tm)^{(d+1)k}) \implies \delta = O(dk \log m). \quad \square$$

² For $x > 1$ if $x / \ln x \leq u$ then $x \leq 2u \ln u$. Hence, if $tm / \log tm \leq dkm$, then $tm = O(dkm \log dkm)$.

Theorem 5.2 Let $(\mathbb{X}_m^d, \mathcal{R}_{dF,k})$ be the range space with $\mathcal{R}_{dF,k}$ being the set of all balls under the discrete Fréchet distance centered at polygonal curves in \mathbb{X}_k^d . The VC dimension is $O(dk \log dk m)$. The shattering dimension is $O(dk \log m)$.

Proof Let $\{S_1, \dots, S_t\} \subseteq \mathbb{X}_m^d$ and $S = \bigcup_i S_i$. Any intersection of a discrete Fréchet ball with $\{S_1, \dots, S_t\}$ is uniquely defined by a sequence $B_1 \cap S, \dots, B_k \cap S$, where B_1, \dots, B_k are balls in \mathbb{R}^d . The number of such sequences can be bounded by $O((tm)^{(d+1)k})$ as in the proof of Theorem 5.1. Enforcing that a sequence contains a valid alignment path only reduces the number of possible distinct sets formed by t curves, and it can be determined using these intersections and the two orderings of B_1, \dots, B_k and of vertices within some $S_j \in \mathbb{X}_m^d$. \square

6 Preliminaries

In this section, we provide a more advanced set of geometric primitives and other known technical results about the VC dimension. We also derive some simple corollaries. Additionally, we provide some basic results about the distances which will couple with the geometric primitives in our proofs for continuous distance measures.

We again consider a ground set $\mathbb{X}_m^d = (\mathbb{R}^d)^m$ which we treat as a set of polygonal curves with points in \mathbb{R}^d of size m . Given such a curve $s \in \mathbb{X}_m^d$, let $V(s)$ be its ordered set of vertices and $E(s)$ its ordered set of edges.

6.1 A Simple Model of Computation

We consider a model of computation that will be useful for modeling primitive geometric sets, and in turn bounding the VC dimension of an associated range space. These will be useful in that they allow the invocation of powerful and general tools to describe range spaces defined by distances between curves. We allow the following operations, which we call *simple operations*:

- the arithmetic operations $+$, $-$, \times , and $/$ on real numbers,
- jumps conditioned on $>$, \geq , $<$, \leq , $=$, and \neq comparisons of real numbers, and
- output 0 or 1.

We say a function requires t simple operations if it can be computed with a circuit of depth t composed only of these simple operations. Note that with the above simple operations, we can also perform logical operations. Furthermore, the lack of a square-root operator creates some challenges when dealing with non-linear geometric objects. Therefore, we prove the following technical lemma showing that we can compare certain expressions involving square roots without computing them explicitly, i.e., only simple operations are needed for the comparison.

Lemma 6.1 Consider values $\alpha, \beta, \gamma, \delta \in \mathbb{R}$ with $\beta, \delta \geq 0$. We can compute the truth value of $\alpha + \sqrt{\beta} \leq \gamma + \sqrt{\delta}$ and $\alpha + \sqrt{\beta} \geq \gamma + \sqrt{\delta}$ using $O(1)$ simple operations.

Proof It suffices to prove the case of $\alpha + \sqrt{\beta} \leq \gamma + \sqrt{\delta}$, as $\alpha + \sqrt{\beta} \geq \gamma + \sqrt{\delta}$ is analogous. We simply show that this comparison is equivalent to a comparison involving only a constant number of simple operations starting from the values $\alpha, \beta, \gamma, \delta$. If

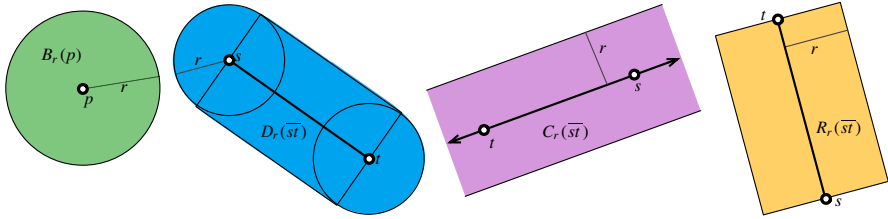


Fig. 1 Illustration of basic shapes in \mathbb{R}^2 , from left to right: a ball $B_r(p)$, a stadium $D_r(\overline{st})$, a cylinder $C_r(\overline{st})$, and a capped cylinder $R_r(\overline{st})$

$\alpha = \gamma$, then $\sqrt{\beta} \leq \sqrt{\delta}$ is equivalent to $\beta \leq \delta$ and we are done. Assuming $\alpha < \gamma$, we get

$$\begin{aligned} \alpha + \sqrt{\beta} \leq \gamma + \sqrt{\delta} &\iff \sqrt{\beta} \leq (\gamma - \alpha) + \sqrt{\delta} \\ &\iff \beta \leq (\gamma - \alpha)^2 + 2(\gamma - \alpha)\sqrt{\delta} + \delta \\ &\iff \beta - (\gamma - \alpha)^2 - \delta \leq 2(\gamma - \alpha)\sqrt{\delta}. \end{aligned}$$

The second equivalence holds because both sides are at least 0. Now, note that the right side of the last inequality is at least 0 and thus, if the left side is negative (which we can check using $O(1)$ simple operations), we are done. Thus, assume the left side is at least 0. Then we can square both sides and obtain a comparison involving only simple operations. Now, if $\gamma < \alpha$, we can do an analogous calculation, where we subtract γ instead of α in the first equivalence. As testing $\gamma < \alpha$ is a simple operation, we can determine which case we are in. \square

6.2 Geometric Primitives

For any $p \in \mathbb{R}^d$ we denote by $B_r(p)$ the ball of radius r , centered at p . For any two points $s, t \in \mathbb{R}^d$, we denote by \overline{st} the line segment from s to t . Whenever we store such a line segment, for technicalities within the lemma below, we store the coordinates of its endpoints s and t . For any two points $s, t \in \mathbb{R}^d$, we define the stadium centered at \overline{st} , $D_r(\overline{st}) = \{x \in \mathbb{R}^d \mid \exists p \in \overline{st} : \|p - x\| \leq r\}$. For any two points $s, t \in \mathbb{R}^d$, we define a cylinder

$$C_r(\overline{st}) = \{x \in \mathbb{R}^d \mid \exists p \in \ell(\overline{st}) : \|p - x\| \leq r\},$$

where $\ell(\overline{st})$ denotes the line supporting the edge \overline{st} . Finally, for any two points $s, t \in \mathbb{R}^d$, we define the capped cylinder centered at \overline{st} : $R_r(\overline{st}) = \{p + u \mid p \in \overline{st} \text{ and } u \in \mathbb{R}^d \text{ s.t. } \|u\| \leq r \text{ and } \langle t - s, u \rangle = 0\}$ (Fig. 1).

For each of these geometric sets, we can determine if a point $x \in \mathbb{R}^d$ is in the set with a constant number of operations under a simple model of computation.

Lemma 6.2 For a point $x \in \mathbb{R}^d$, and any set of the form $B_r(p)$, $D_r(\overline{st})$, $C_r(\overline{st})$, or $R_r(\overline{st})$, we can determine if x is in that set (returns 1, otherwise 0) using $O(d)$ simple operations.

Proof For the ball $B_r(p)$ we can compute a distance $\|x - p\|^2$ in $O(d)$ time, and determine inclusion with a comparison to r^2 . For the cylinder $C_r(\overline{st})$ we can compute the closest point to x on this line as

$$\pi_{\overline{st}}(x) = t + \frac{(s - t)\langle s - t, x - t \rangle}{\|s - t\|^2}.$$

Then we can determine inclusion by comparing $\|\pi_{\overline{st}}(x) - x\|^2$ to r^2 . For the capped cylinder $R_r(\overline{st})$ we also need to compare $\|\pi_{\overline{st}}(x) - t\|^2$ and $\|\pi_{\overline{st}}(x) - s\|^2$ to see if either of these terms is greater than $\|s - t\|^2$. For the stadium $D_r(\overline{st})$ we determine inclusion if x is in any of $R_r(\overline{st})$, $B_r(s)$, or $B_r(t)$. □

6.3 Bounding the VC Dimension

For range spaces defined on continuous curves, our proofs use a powerful theorem from Goldberg and Jerrum [27] as improved and restated by Anthony and Bartlett [7]. It allows one to easily bound the VC dimension of geometric range spaces under our simple model of computation.

Theorem 6.3 ([7, Thm. 8.4]) Suppose h is a function from $\mathbb{R}^d \times \mathbb{R}^n$ to $\{0, 1\}$ and let

$$H = \{x \mapsto h(\alpha, x) \mid \alpha \in \mathbb{R}^d\}$$

be the range set determined by h with preimage of 1. Suppose that h can be computed by an algorithm that takes as input the pair $(\alpha, x) \in \mathbb{R}^d \times \mathbb{R}^n$ and returns $h(\alpha, x)$ after no more than t simple operations. Then, the VC dimension of H is $\leq 4d(t + 2)$.

An example implication can be seen for geometric sets via Lemma 6.2. Note that this implies any VC dimension upper bound proven in this approach applies to both the range space and its dual range space because the function h is unchanged and the ranges can still be described by $O(d)$ real coordinates.

Corollary 6.4 For range spaces defined on \mathbb{R}^d with geometric sets $B_r(p)$, $D_r(\overline{st})$, $C_r(\overline{st})$, or $R_r(\overline{st})$ as ranges, the VC dimension is $O(d^2)$. The same $O(d^2)$ VC dimension bound holds for the corresponding dual range spaces, with ground sets as the geometric sets, and ranges defined by stabbing using points in \mathbb{R}^d .

Note that these bounds are not always tight. Specifically, because the VC dimension for ranges defined geometrically by balls $B_r(p)$ is $O(d)$ [30]. Moreover, the VC dimension of range spaces defined by cylinders $C_r(\overline{st})$ is known to be $O(d)$ [4]. The ranges defined by capped cylinders $R_r(\overline{st})$ are the intersection of a cylinder and two halfspaces, each with VC dimension $O(d)$ and hence, by the composition theorem [12], this full range space also has VC dimension $O(d)$. Finally, the stadium $D_r(\overline{st})$ is defined by the

union of a capped cylinder $R_r(\overline{st})$ and two balls $B_r(s)$ and $B_r(t)$; hence, again by the composition theorem [12], its VC dimension is $O(d)$.

However, it is not clear that these improved bounds hold for the dual range spaces, aside from the case of B_r . Moreover, when the ground set X of the range space (X, \mathcal{R}) is not \mathbb{R}^d , then we need to be cautious in using the k -fold composition theorem [12], which bounds the VC dimension of complex range spaces derived as the logical intersection or union of simpler range spaces with bounded VC dimension. In the case of a ground set $X = \mathbb{R}^d$, logical and geometric intersections are the same, but for other ground sets (like dual objects, or line segments \mathbb{X}_2^d) this is not necessarily the case. For instance, a line segment $e \in \mathbb{X}_2^d$ may intersect a ball B_r and also a halfspace H while not intersecting the intersection $B_r \cap H$.

6.4 Representation by Predicates

In order to prove bounds on the VC dimension of range spaces defined on continuous curves, we establish sets of geometric predicates which are sufficient to determine if two curves have distance at most r to each other. Analyzing the range spaces associated with these predicates (over all possible radii r) allows us to compose them further and to establish VC dimension bounds for the range space induced by the corresponding distance measure. For the Fréchet and weak Fréchet distance, the predicates mirror those used in range searching data structures [1,2]. And for the Hausdorff distance on continuous curves, the predicates are derived from the Voronoi diagram [5]. The technical challenges for each case are similar, but require different analyses.

7 The Hausdorff Distance

We consider the range space $(\mathbb{X}_m^d, \mathcal{R}_{H_k}^r)$, where $\mathcal{R}_{H_k}^r$ denotes the set of all balls, of radius r , centered at curves in \mathbb{X}_k^d , under the Hausdorff distance.³ We also consider the same problems under both directed versions of the Hausdorff distance, and their induced range spaces $(\mathbb{X}_m^d, \mathcal{R}_{\overrightarrow{H}_k}^r)$ and $(\mathbb{X}_m^d, \mathcal{R}_{\overleftarrow{H}_k}^r)$, where $\mathcal{R}_{\overrightarrow{H}_k}^r$ denotes the set of all balls of radius r under the directed Hausdorff distance from curves in \mathbb{X}_k^d , and $\mathcal{R}_{\overleftarrow{H}_k}^r$ denotes the set of all balls of radius r under the directed Hausdorff distance from curves to \mathbb{X}_k^d .

7.1 Hausdorff Distance Predicates

Consider two sets of line segments A and B such that any two segments that belong to the same set have disjoint interiors. Consider the Voronoi diagram of the vertices and open segments of B : each element of B (i.e., open segment or vertex) is assigned to a Voronoi cell which is the set of points that are closer to this element than to any other element (see Fig. 2). According to Alt et al. [5], the critical points for the

³ The proofs in this section are written for polygonal curves in \mathbb{X}_m^d , but they readily extend to (not necessarily connected) sets of line segments in \mathbb{R}^d of cardinality $m' = (m - 1)/2$.

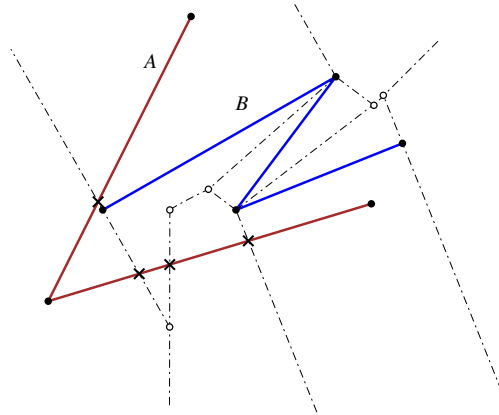


Fig. 2 Two polygonal lines A and B . The critical points for directed Hausdorff distance $d_{\vec{H}}(A, B)$ are either at some vertex of A or at some intersection point of A with the boundary of a Voronoi cell of B

directed Hausdorff distance $d_{\vec{H}}(A, B)$ occur either at some vertex of A or at some intersection point of A with the boundary of a Voronoi cell of B . Thus, we need a predicate for encoding the first type of event where the distance is assumed at a vertex of A . Additionally, we need a predicate for testing if a line supporting an edge intersects the intersection of two stadiums; see Fig. 3 for an illustration in \mathbb{R}^2 .

Consider any two polygonal curves $s \in \mathbb{X}_m^d$ and $q \in \mathbb{X}_k^d$. In order to encode the intersection of polygonal curves with metric balls under the Hausdorff metric, we will first define a subset of \mathbb{R}^d , a *double-stadium*, defined by two line segments $\{e_1, e_2\}$ and a radius r as

$$D_{r,2}(e_1, e_2) = D_r(e_1) \cap D_r(e_2).$$

We use the notation $\overline{uv} \in D_{r,2}(e_1, e_2)$ to indicate that the line $\ell(\overline{uv})$ which supports \overline{uv} intersects with the double-stadium, i.e., it fulfills

$$\ell(\overline{uv}) \cap D_{r,2}(e_1, e_2) \neq \emptyset.$$

We will make use of the following predicates:

- P_1 (*Vertex-edge (horizontal)*) Given an edge of s , $\overline{s_j s_{j+1}}$, and a vertex q_i of q , this predicate returns true iff there exists a point $p \in \overline{s_j s_{j+1}}$, such that $\|p - q_i\| \leq r$.
- P_2 (*Vertex-edge (vertical)*) Given an edge of q , $\overline{q_i q_{i+1}}$, and a vertex s_j of s , this predicate returns true iff there exists a point $p \in \overline{q_i q_{i+1}}$, such that $\|p - s_j\| \leq r$.
- P_3 (*d-stadium-line (horizontal)*) Given an edge of q , $\overline{q_i q_{i+1}}$, and two edges of s , $\{e_1, e_2\} \subset E(s)$, this predicate is equal to $\overline{q_i q_{i+1}} \in D_{r,2}(e_1, e_2)$.
- P_4 (*d-stadium-line (vertical)*) Given one edge of s , $\overline{s_j s_{j+1}}$, and two edges of q , $\{e_1, e_2\} \subset E(q)$, this predicate is equal to $\overline{s_j s_{j+1}} \in D_{r,2}(e_1, e_2)$.

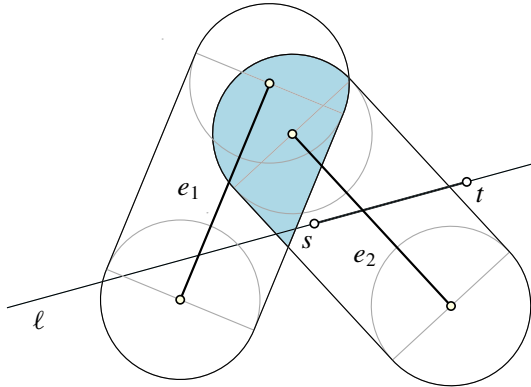


Fig. 3 Illustration of the predicate P_3 in \mathbb{R}^2 : The predicate evaluates to true if and only if the triple intersection of the line ℓ supporting $\overline{s\bar{t}}$ with the two stadiums centered at edges e_1 and e_2 is non-empty. Note that $\overline{s\bar{t}}$ may lie outside of the intersection

Lemma 7.1 For any two polygonal curves s, q , given the truth values of the predicates P_1, P_3 one can determine whether $d_{\overline{H}}(q, s) \leq r$. Similarly, given the truth values of the predicates P_2, P_4 one can determine whether $d_{\overline{H}}(s, q) \leq r$.

Proof We first assume for the sake of simplicity that q is a line segment with endpoints q_1 and q_2 . We claim that $d_{\overline{H}}(q, s) \leq r$ if and only if there exists a sequence of edges $\overline{s_{j_1} s_{j_1+1}}, \overline{s_{j_2} s_{j_2+1}}, \dots, \overline{s_{j_v} s_{j_v+1}}$ for some integer value v , such that the predicates $P_1(q_1, \overline{s_{j_1} s_{j_1+1}}), P_1(q_2, \overline{s_{j_v} s_{j_v+1}})$ both evaluate to true and the conjugate

$$\bigwedge_{i=1}^{v-1} P_3(\overline{q_1 q_2}, \overline{s_{j_i} s_{j_i+1}}, \overline{s_{j_{i+1}} s_{j_{i+1}+1}})$$

evaluates to true. Assume such a sequence of edges exists. In this case, there exists a sequence of points p_1, \dots, p_v on the line supporting q , with $p_1 = q_1, p_v = q_2$, and such that for $1 \leq i < v, p_i, p_{i+1} \in D_r(\overline{s_{j_i} s_{j_i+1}})$. That is, two consecutive points of the sequence are contained in the same stadium. Indeed, for $i = 1$ we have $p_1 = q_1$ and $p_2 \in \overline{s_{j_1} s_{j_1+1}}$ since the corresponding P_1 and P_3 predicates evaluate to true:

$$P_1(q_1, \overline{s_{j_1} s_{j_1+1}}), \quad P_3(\overline{q_1 q_2}, \overline{s_{j_1} s_{j_1+1}}, \overline{s_{j_2} s_{j_2+1}}).$$

Likewise, for $i = v - 1$, it is implied by the corresponding predicates $P_1(q_2, \overline{s_{j_v} s_{j_v+1}})$ and $P_3(\overline{q_1 q_2}, \overline{s_{j_{v-1}} s_{j_{v-1}+1}}, \overline{s_{j_v} s_{j_v+1}})$. For the remaining $1 < i < v - 1$, it follows from the conditions given by the specified P_3 predicates. Now, since each stadium is a convex set, it follows that each line segment connecting two consecutive points of this sequence p_i, p_{i+1} is contained in one of the stadiums. Note that the set of line segments obtained this way forms a connected polygonal curve which fully covers the

line segment q . It follows that

$$q \subseteq \bigcup_{0 \leq i < v} \overline{p_i p_{i+1}} \subseteq \bigcup_{0 \leq i < v} D_r(\overline{s_{j_i} s_{j_{i+1}}}) \subseteq \bigcup_{e \in E(s)} D_r(e).$$

Therefore, any point on q is within distance r of some point on s and thus $d_{\vec{H}}(q, s) \leq r$.

Now, for the other direction of the proof, assume that $d_{\vec{H}}(q, s) \leq r$. The definition of the directed Hausdorff distance implies that

$$q \subseteq \bigcup_{e \in E(s)} D_r(e),$$

since any point on the line segment q must be within distance r of some point on the curve s . Consider the intersections of the line segment q with the boundaries of stadiums

$$q \cap \bigcup_{e \in E(s)} \partial D_r(e).$$

Let w be the number of intersection points and let $v = w + 2$. We claim that this implies that there exists a sequence of edges $\overline{s_{j_1} s_{j_1+1}}, \overline{s_{j_2} s_{j_2+1}}, \dots, \overline{s_{j_v} s_{j_v+1}}$ with the properties stated above. Let $p_1 = q_1, p_v = q_2$, and let p_i for $1 < i < v$ be the intersection points ordered in the direction of the line segment q . By construction, it must be that each p_i for $1 < i < v$ is contained in the intersection of two stadiums, since it is the intersection with the boundary of a stadium and the entire edge is covered by the union of stadiums. Moreover, two consecutive points p_i, p_{i+1} are contained in exactly the same subset of stadiums—otherwise there would be another intersection point with the boundary of a stadium in between p_i and p_{i+1} . This implies a set of true predicates of type P_3 with the properties defined above. The predicates of type P_1 follow trivially from the definition of the directed Hausdorff distance. This concludes the proof of the other direction.

In general, for any polygonal curve $q \in \mathbb{X}_k^d$ with vertices q_1, \dots, q_k , we have that

$$d_{\vec{H}}(q, s) \leq r \iff \bigwedge_{i=1}^{k-1} [d_{\vec{H}}(\overline{q_i q_{i+1}}, s) \leq r].$$

Thus, we can apply the arguments above to each edge of q individually. Similarly, we can prove that given the truth values of the predicates P_2, P_4 one can determine whether $d_{\vec{H}}(s, q) \leq r$, by an argument symmetric to the above. □

7.2 Hausdorff Distance VC Dimension Bound

Now, we want to show that we can compute a representation of the interval of intersection of a line and a capped cylinder using only $O(d)$ simple operations. This

representation then allows us to compare such intervals using Lemma 6.1. The appropriate ground set is over two points $q_j, q_t \in \mathbb{R}^d$, where for notational simplicity we reuse \mathbb{X}_2^d . Furthermore, for each segment $\overline{st} \in \mathbb{X}_2^d$, recall that $\ell(\overline{st})$ is the line that supports it.

Lemma 7.2 *Given a line $\ell(\overline{st})$ with $\overline{st} \in \mathbb{X}_2^d$ and a capped cylinder $R_r(\overline{uv})$ with $\overline{uv} \in \mathbb{X}_2^d$, the intersection $\ell(\overline{st}) \cap R_r(\overline{uv})$ of those two objects is either*

$$\{s + (t - s)x \mid x \in [\alpha + \sqrt{\beta}, \gamma + \sqrt{\delta}] \subseteq \mathbb{R}\},$$

where $\alpha, \beta, \gamma, \delta \in \mathbb{R}$ can be computed using $O(d)$ simple operations, or it is empty.

Proof We first compute the intersection of the infinite cylinder $C_r(\overline{uv})$ with the line $\ell(\overline{st})$. Let $f(x) = u - (v - u)x$ be the line $\ell(\overline{uv})$ parametrized by $x \in \mathbb{R}$ and $g(y) = s + (t - s)y$ the line $\ell(\overline{st})$ parametrized by $y \in \mathbb{R}$. We describe all values x, y parameterizing points in this intersection by quantifying the boundaries of this set. All points in the intersection of $\ell(\overline{st})$ with the boundary of the infinite cylinder $C_r(\overline{uv})$ are described by

$$\begin{aligned} \|g(y) - f(x)\|^2 = r^2 &\iff \|s + (t - s)y - u - (v - u)x\|^2 = r^2 \\ &\iff \sum_{i=1}^d (s_i + (t_i - s_i)y - u_i - (v_i - u_i)x)^2 = r^2. \end{aligned}$$

Let $z_i(y) = s_i + (t_i - s_i)y - u_i$. We obtain

$$\begin{aligned} \sum_{i=1}^d (z_i(y) - (v_i - u_i)x)^2 &= r^2 \\ \iff \sum_{i=1}^d ((v_i - u_i)^2 x^2 - 2z_i(y)(v_i - u_i)x + z_i(y)^2) - r^2 &= 0. \end{aligned}$$

For any fixed y , this is a quadratic equation in x and the discriminant is

$$h(y) = \left(\sum_{i=1}^d 2z_i(y)(v_i - u_i) \right)^2 - 4 \sum_{i=1}^d (v_i - u_i)^2 \cdot \left(\sum_{i=1}^d z_i(y)^2 - r^2 \right).$$

Note that the quadratic equation has one solution exactly for those points on $\ell(\overline{st})$ which have distance r from $\ell(\overline{uv})$, because the ball around those points intersects $\ell(\overline{uv})$ exactly once. Those are also the points which define the boundary of $\ell(\overline{st}) \cap R_r(\overline{uv})$. Thus, we want to solve $h(y) = 0$. As $z_i(y)$ is linear in y , we obtain a quadratic equation in y . Note that all coefficients of the quadratic equation can be computed in $O(d)$ simple operations. Both solutions of this equation are of the form $\alpha \pm \sqrt{\beta}$. If $\beta < 0$, then the intersection is empty. Otherwise, we obtain an intersection interval $[\alpha - \sqrt{\beta}, \alpha + \sqrt{\beta}]$ for the infinite cylinder.

To obtain the intersection with the capped cylinder, we first compute the intersection of $\ell(\overline{st})$ with the top and bottom hyperplanes of the cylinder. The two planes are given by all $p \in \mathbb{R}^d$ which satisfy $(p - u)(v - u) = 0$ and $(p - v)(v - u) = 0$, respectively. By plugging the line equation into the hyperplane formulas, we get the intersection points. For the first plane we thereby obtain

$$(g(y) - u)(v - u) = 0 \iff (s + (t - s)y - u)(v - u) = 0$$

which resolves to

$$y = -\frac{(s - u)(v - u)}{(t - s)(v - u)}.$$

The intersection with the second plane is analogous. Thus, we again obtain an interval for y such that the values in this interval induce the intersection points between the planes. Again $O(d)$ simple operations are sufficient to compute the boundaries of this interval.

To obtain the intersection with the capped cylinder (not just with its boundary planes), we intersect the two intervals we obtained for the intersection with the infinite cylinder as well as the boundary planes of the capped cylinder. As computing the intersection of intervals is simply taking the minimum/maximum, we can use Lemma 6.1 to do this in $O(1)$ simple operations. The values for $\alpha, \beta, \gamma, \delta$ are then given by the intersection interval boundaries which are chosen from the boundaries of the intersection interval of the planes of the capped cylinder and the infinite cylinder. \square

Additionally, the following lemma holds, which states that we can express an intersection of a ball and a line with an interval of the form as in the previous lemma.

Lemma 7.3 *Given a line $\ell(\overline{st})$ with $\overline{st} \in \mathbb{X}_2^d$ and a ball $B_r(c)$ centered at c , the intersection $\ell(\overline{st}) \cap B_r(c)$ of those two objects is either*

$$\{s + (t - s)x \mid x \in [\alpha + \sqrt{\beta}, \gamma + \sqrt{\delta}] \subseteq \mathbb{R}\},$$

where $\alpha, \beta, \gamma, \delta \in \mathbb{R}$ can be computed using $O(d)$ simple operations, or it is empty.

Proof The intersection is given by the x fulfilling $\|s + (t - s)x - c\|^2 \leq r^2$. To determine the extremal values for x which satisfy this inequality is a quadratic equation in x . Solving it, we obtain an intersection interval as required. \square

Having proven those technical lemmas, we are now ready to start our argument for bounding the VC dimension. We argue that the truth values for predicate P_1 over all possible inputs are uniquely defined by the set

$$- P_1^r(q, s) = \{D_r(\overline{s_i s_{i+1}}) \cap V(q) \mid \overline{s_i s_{i+1}} \in E(s)\}.$$

Similarly, the truth values for predicate P_2 are uniquely defined by the set

$$- P_2^r(q, s) = \{D_r(\overline{q_i q_{i+1}}) \cap V(s) \mid \overline{q_i q_{i+1}} \in E(q)\}.$$

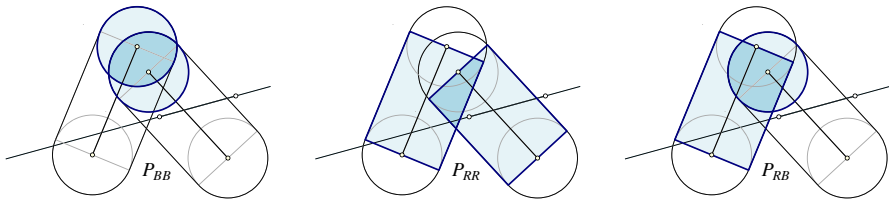


Fig. 4 Illustration in \mathbb{R}^2 of predicates used in the proof of Lemma 7.4 for the example given in Fig. 3

Then the predicates P_3 and P_4 induce sets (where effectively $P_4(q, s) = P_3(s, q)$)

- $P_3^r(q, s) = \{(e_1, e_2, e_3) \in E(s) \times E(s) \times E(q) \mid e_3 \in D_{r,2}(e_1, e_2)\}$,
- $P_4^r(q, s) = \{(e_1, e_2, e_3) \in E(q) \times E(q) \times E(s) \mid e_3 \in D_{r,2}(e_1, e_2)\}$.

We require a technical proof, bounding the VC dimension of the range space defined on segments with ranges defined by double-stadiums. To this end, let

$$\mathcal{D}_2^d = \{\{\overline{st} \in \mathbb{X}_2^d \mid \overline{st} \in D_{r,2}(e_1, e_2)\} \mid e_1, e_2 \in \mathbb{X}_2^d, r > 0\}$$

be the families of subsets of line segments $\overline{st} \in \mathbb{X}_2^d$ whose supported lines $\ell(\overline{st})$ intersect a common double-stadium $D_{r,2}(e_1, e_2)$. We are now ready to state and prove the following lemma.

Lemma 7.4 *The VC dimension of the range space $(\mathbb{X}_2^d, \mathcal{D}_2^d)$ and of the associated dual range space is $O(d^2)$.*

Proof The predicate which determines whether a line ℓ intersects a double-stadium $D_{r,2}(e_1, e_2)$ can be implemented by taking the logical-or over $O(1)$ calls to the following predicates (see Fig. 4 for an illustration):

- P_{BB} : checks whether ℓ intersects $D_{r,2}(e_1, e_2)$ in the intersection of two radius r balls,
- P_{RR} : checks whether ℓ intersects $D_{r,2}(e_1, e_2)$ in the intersection of two radius r capped cylinders,
- P_{RB} : checks whether ℓ intersects $D_{r,2}(e_1, e_2)$ in the intersection of one ball and one capped cylinder, both of radius r .

For all predicates we first compute the intersection interval of the capped cylinder or ball using Lemmas 7.2 or 7.3. Applying Lemma 6.1, we can then compute the intersection of these two intersection intervals by comparing their bounds, obtaining an interval of the form $[\alpha + \sqrt{\beta}, \gamma + \sqrt{\delta}]$. Again using Lemma 6.1, we test if $\alpha + \sqrt{\beta} \leq \gamma + \sqrt{\delta}$, thereby checking if the intersection is non-empty. Thus, all three of the above predicates can be computed in $O(d)$ simple operations. Because each predicate requires $O(d)$ simple operations, and we need to perform a logical-or over $O(1)$ of these predicates, it implies range inclusion $e \in D_{r,2}$ and can be determined with $O(d)$ simple operations. Hence by Theorem 6.3 the VC dimension is $O(d^2)$. Since an element of the dual range space is also defined by $O(d)$ real values, and the same operations can be applied, the dual range space also has VC dimension $O(d^2)$. \square

Using the above lemmas, we now get the following theorems.

Theorem 7.5 Let $\vec{\mathcal{R}}_{H,k}$ be the set of all balls, under the directed Hausdorff distance from polygonal curves in \mathbb{X}_k^d . The VC dimension of $(\mathbb{X}_m^d, \vec{\mathcal{R}}_{H,k})$ is $O(d^2k \log dkm)$. The shattering dimension of $(\mathbb{X}_m^d, \vec{\mathcal{R}}_{H,k})$ is $O(d^2k \log m)$.

Proof Let $S \subset \mathbb{X}_m^d$ be a set of t polygonal curves and let $q \in \mathbb{X}_k^d$. By Lemma 7.1, the set $\{s \in S \mid d_{\vec{H}}(q, s) \leq r\}$ is uniquely defined by the sets

$$\bigcup_{s \in S} P_1^r(q, s), \quad \bigcup_{s \in S} P_3^r(q, s).$$

The number of all possible sets $\bigcup_{r \geq 0} \bigcup_{s \in S} P_1^r(q, s)$ is bounded by $(tm)^{O(d^2k)}$. This follows by the upper bound of Corollary 6.4, on the VC dimension of the range space having as ground set the set of stadiums and ranges corresponding to stabbing points, and the fact that we need to consider k vertices for the query curve. Furthermore, by Lemma 7.4, we are able to bound the number of all possible sets $\bigcup_{r \geq 0} \bigcup_{s \in S} P_3^r(q, s)$ as $(tm)^{O(d^2k)}$. The k term in the exponent arises because we consider all k edges of q for predicate P_3 . Hence,

$$2^t \leq (tm)^{O(d^2k)} \implies t = O(d^2k \log dkm).$$

We can similarly bound the shattering dimension δ ,

$$t^\delta \leq (tm)^{O(d^2k)} \implies \delta = O(d^2k \log m). \quad \square$$

Theorem 7.6 Let $\overleftarrow{\mathcal{R}}_{H,k}$ be the set of all balls, under the directed Hausdorff distance to polygonal curves in \mathbb{X}_k^d . The VC dimension of $(\mathbb{X}_m^d, \overleftarrow{\mathcal{R}}_{H,k})$ is $O(d^2k^2 \log dkm)$. The shattering dimension of $(\mathbb{X}_m^d, \overleftarrow{\mathcal{R}}_{H,k})$ is $O(d^2k^2 \log m)$.

Proof Let $S \subset \mathbb{X}_m^d$ be a set of t polygonal curves and let $q \in \mathbb{X}_k^d$. By Lemma 7.1, the set $\{s \in S \mid d_{\overleftarrow{H}}(q, s) \leq r\}$ is uniquely defined by the sets

$$\bigcup_{s \in S} P_2^r(q, s), \quad \bigcup_{s \in S} P_4^r(q, s).$$

The number of all possible sets $\bigcup_{r \geq 0} \bigcup_{s \in S} P_2^r(q, s)$ is bounded by $(tm)^{O(d^2k)}$. This follows by the upper bound of Corollary 6.4, on the VC dimension of range spaces with points as the ground set and stadiums as ranges, and the fact that we need to consider one stadium for each of the $k - 1$ query edges. Furthermore, by Lemma 7.4, we are able to bound the number of all possible sets $\bigcup_{r \geq 0} \bigcup_{s \in S} P_4^r(q, s)$ as $(tm)^{O(d^2k^2)}$. The k^2 term in the exponent arises because we consider $\Theta(k^2)$ pairs of edges of q for predicate P_4 . Now,

$$2^t \leq (tm)^{O(d^2k^2)} \implies t = O(d^2k^2 \log dkm).$$

We can similarly bound the shattering dimension δ ,

$$t^\delta \leq (tm)^{O(d^2k^2)} \implies \delta = O(d^2k^2 \log m). \quad \square$$

Theorem 7.7 *Let $\mathcal{R}_{H,k}$ be the set of all balls, under the symmetric Hausdorff distance in \mathbb{X}_k^d . The VC dimension of $(\mathbb{X}_m^d, \mathcal{R}_{H,k})$ is $O(d^2k^2 \log dkm)$. The shattering dimension of $(\mathbb{X}_m^d, \mathcal{R}_{H,k})$ is $O(d^2k^2 \log m)$.*

Proof Lemma 7.1 implies that the set $\{s \in S \mid d_H(q, s) \leq r\}$ is uniquely defined by the sets

$$\bigcup_{s \in S} P_1^r(q, s), \quad \bigcup_{s \in S} P_2^r(q, s), \quad \bigcup_{s \in S} P_3^r(q, s), \quad \bigcup_{s \in S} P_4^r(q, s).$$

Now bounding the number of all possible such sets, as we did in the proofs of Theorems 7.5 and 7.6, implies the statement. \square

8 The Fréchet Distance

We consider the range spaces $(\mathbb{X}_m^d, \mathcal{R}_{F_k})$ and $(\mathbb{X}_m^d, \mathcal{R}_{wF_k})$, where \mathcal{R}_{F_k} (resp. \mathcal{R}_{wF_k}) denotes the set of all balls, centered at curves in \mathbb{X}_k^d , under the Fréchet (resp. weak Fréchet) distance.

8.1 Fréchet Distance Predicates

It is known that the Fréchet distance between two polygonal curves can be attained, either at a distance between their endpoints, at a distance between a vertex and a line supporting an edge, or at the common distance of two vertices with a line supporting an edge. The third type of event is sometimes called monotonicity event, since it happens when the weak Fréchet distance is smaller than the Fréchet distance. In this sense, our representation of the ball of radius r under the Fréchet distance is based on the following predicates, some of which we already used in the last section. Let $s \in \mathbb{X}_m^d$ with vertices s_1, \dots, s_m and $q \in \mathbb{X}_k^d$ with vertices q_1, \dots, q_k .

P_1 (Vertex-edge (horizontal)) As defined in Sect. 7.

P_2 (Vertex-edge (vertical)) As defined in Sect. 7.

P_5 (Endpoints (start)) This predicate returns true if and only if $\|s_1 - q_1\| \leq r$.

P_6 (Endpoints (end)) This predicate returns true if and only if $\|s_m - q_k\| \leq r$.

P_7 (Monotonicity (horizontal)) Given two vertices of s , s_j and s_t with $j < t$, and an edge of q , $\overline{q_i q_{i+1}}$, this predicate returns true if there exist two points p_1 and p_2 on the line supporting the directed edge, such that p_1 appears before p_2 on this line, and such that $\|p_1 - s_j\| \leq r$ and $\|p_2 - s_t\| \leq r$.

P_8 (Monotonicity (vertical)) Given two vertices of q , q_i and q_t with $i < t$, and a directed edge of s , $\overline{s_j s_{j+1}}$, this predicate returns true if there exist two points p_1

and p_2 on the line supporting the directed edge, such that p_1 appears before p_2 on this line, and such that $\|p_1 - q_i\| \leq r$ and $\|p_2 - q_t\| \leq r$.

Predicate P_8 is illustrated in Fig. 5. Predicate P_7 is symmetric.

Lemma 8.1 ([1, Lem. 9]) *Given the truth values of all predicates $P_1, P_2, P_5, P_6, P_7, P_8$ of two curves s and q for a fixed value of r , one can determine if $d_F(s, q) \leq r$.*

Predicates P_1, P_2, P_5, P_6 are sufficient for representing metric balls under the weak Fréchet distance. We include a proof for the sake of completeness.

Lemma 8.2 *Given the truth values of all predicates P_1, P_2, P_5, P_6 of two curves s and q for a fixed value of r , one can determine if $d_{wF}(s, q) \leq r$.*

Proof Alt and Godau [6] describe an algorithm for computing the weak Fréchet distance which can be used here. In particular, one can construct an edge-weighted grid graph on the cells (edge–edge pairs) of the parametric space of the two polygonal curves, and subsequently compute a bottleneck-shortest path from the pair of first edges to the pair of last edges along the two curves. We can use edge weights in $\{0, 1\}$ to encode if the corresponding vertex-edge pair has distance at most r , as given by the predicates P_1 and P_2 . If and only if there exists a bottleneck shortest path of cost 0 and the endpoint conditions are satisfied (as given by the predicates P_5 and P_6), the weak Fréchet distance between q and s is at most r . \square

8.2 Fréchet Distance VC Dimension Bounds

We first consider the range space $(\mathbb{X}_m^d, \mathcal{R}_{wF,k})$, where $\mathcal{R}_{wF,k}$ is the set of all balls under the weak Fréchet distance centered at curves in \mathbb{X}_k^d . The main task is to translate the predicates P_1, P_2, P_5, P_6 into simple range spaces, and then bound their associated VC dimensions. Consider any two polygonal curves $s \in \mathbb{X}_m^d$ and $q \in \mathbb{X}_k^d$. In order to encode the intersection of polygonal curves with metric balls, we will make use of the sets $P_1^r(q, s), P_2^r(q, s)$, which are defined in Sect. 7, and the following sets:

- $P_5^r(q, s) = B_r(q_1) \cap V(s)$,
- $P_6^r(q, s) = B_r(q_k) \cap V(s)$,

Theorem 8.3 *Let $\mathcal{R}_{wF,k}$ be the set of balls under the weak Fréchet metric centered at polygonal curves in \mathbb{X}_k^d . The VC dimension of $(\mathbb{X}_m^d, \mathcal{R}_{wF,k})$ is $O(d^2k \log dkm)$. The shattering dimension of $(\mathbb{X}_m^d, \mathcal{R}_{wF,k})$ is $O(d^2k \log m)$.*

Proof If S is a set of t polygonal curves of complexity m , $\{s \in S \mid d_{wF}(s, q) \leq r\}$ is uniquely defined by the sets

$$\bigcup_{s \in S} P_1^r(q, s), \quad \bigcup_{s \in S} P_2^r(q, s), \quad \bigcup_{s \in S} P_5^r(q, s), \quad \bigcup_{s \in S} P_6^r(q, s).$$

The number of all possible sets $\bigcup_{r \geq 0} \bigcup_{s \in S} P_1^r(q, s)$ and the number of all possible sets $\bigcup_{r \geq 0} \bigcup_{s \in S} P_2^r(q, s)$ are both bounded by $(tm)^{O(d^2k)}$ by Corollary 6.4 using set $D_r(\overline{st})$, and by considering the dual range space, respectively.

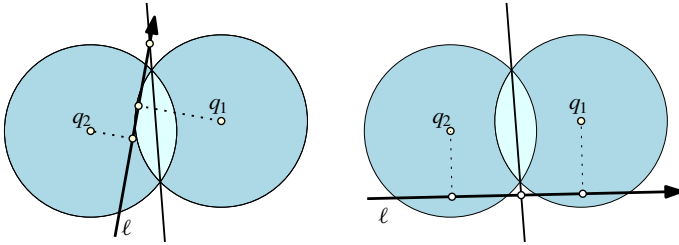


Fig. 5 Illustration of predicate P_8 in \mathbb{R}^2 with line ℓ and the two disks centered at q_1 and q_2 . In these examples, the projection of q_2 onto ℓ appears before the projection of q_1 onto ℓ along the direction of ℓ and the intersection of ℓ with the bisector lies outside of the lens formed by the two disks. On the left, the predicate is satisfied by setting $p_1 = p_2 = \pi_{\overline{q_1q_2}}(q_1)$. On the right, the predicate evaluates to false

Notice that the number of all possible sets $\bigcup_{r \geq 0} \bigcup_{s \in S} P_5^r(q, s)$ is bounded by $(tm)^{O(d)}$. The same holds for the number of all possible sets $\bigcup_{r \geq 0} \bigcup_{s \in S} P_6^r(q, s)$. Hence,

$$2^t \leq (tm)^{O(d^2k)} \implies t = O(d^2k \log dkm).$$

We can similarly bound the shattering dimension δ ,

$$t^\delta \leq (tm)^{O(d^2k)} \implies \delta = O(d^2k \log m). \quad \square$$

We now consider the range space $(\mathbb{X}_m^d, \mathcal{R}_{F,k})$, where $\mathcal{R}_{F,k}$ denotes the set of all balls, centered at curves in \mathbb{X}_k^d , under the Fréchet distance. The approach is the same as with the weak Fréchet distance, except we also need to bound the VC dimension of range spaces associated with predicates P_7 and P_8 to encode monotonicity. For that, we can simply appeal to Theorem 6.3.

We need to define a set to represent predicates P_7 and P_8 . To this end, we again use \mathbb{X}_2^d to represent the set of all segments in \mathbb{R}^d . Given radius $r \geq 0$ and a line segment \overline{st} , we define $M_r(\overline{st})$ to be the set containing all pairs of points (q_1, q_2) for which there exist $p_1, p_2 \in \ell$, where \overline{st} supports ℓ , such that

- $\|p_1 - q_1\| \leq r$ and $\|p_2 - q_2\| \leq r$,
- p_1 is less than p_2 along the line, as $\langle p_1, t - s \rangle \leq \langle p_2, t - s \rangle$.

The predicate P_7 is satisfied if and only if $(s_j, s_t) \in M_r(\overline{q_iq_{i+1}})$ and predicate P_8 is satisfied if and only if $(q_i, q_t) \in M_r(\overline{s_j s_{j+1}})$. Finally, we define $\mathcal{M} = \{M_r(\overline{st}) \mid \overline{st} \in \mathbb{X}_2^d, r \geq 0\}$ to be the set of all relevant ranges.

Corollary 8.4 *The VC dimension of the range space $(\mathbb{X}_2^d, \mathcal{M})$, and of the associated dual range space, is $O(d^2)$.*

Proof The corollary directly follows from Lemma 7.4 by collapsing the stadiums to circles. □

We define sets to correspond with predicates P_7 and P_8 :

- $P_7^r(q, s) = \{(s_j, s_t) \in V(s) \times V(s) \mid (s_j, s_t) \in M_r(\overline{q_i q_{i+1}}) \text{ and } \overline{q_i q_{i+1}} \in E(q)\}.$
- $P_8^r(q, s) = \{(q_i, q_t) \in V(q) \times V(q) \mid (s_i, s_t) \in M_r(\overline{s_j s_{j+1}}) \text{ and } \overline{s_j s_{j+1}} \in E(s)\}.$

Theorem 8.5 *Let $\mathcal{R}_{F,k}$ be the set of all balls, under the Fréchet distance, centered at polygonal curves in \mathbb{X}_k^d . The VC dimension of $(\mathbb{X}_m^d, \mathcal{R}_{F,k})$ is $O(d^2 k^2 \log dkm)$. The shattering dimension of $(\mathbb{X}_m^d, \mathcal{R}_{F,k})$ is $O(d^2 k^2 \log m)$.*

Proof Due to Lemma 8.1, if $S \subset \mathbb{X}_m^d$ is a set of t polygonal curves and $q \in \mathbb{X}_k^d$, the set $\{s \in S \mid d_F(s, q) \leq r\}$ is uniquely defined by the sets

$$\begin{aligned} & \bigcup_{s \in S} P_1^r(q, s), \quad \bigcup_{s \in S} P_2^r(q, s), \quad \bigcup_{s \in S} P_5^r(q, s), \\ & \bigcup_{s \in S} P_6^r(q, s), \quad \bigcup_{s \in S} P_7^r(q, s), \quad \bigcup_{s \in S} P_8^r(q, s). \end{aligned}$$

As in the proof of Theorem 8.3, the number of all possible sets

$$\bigcup_{r \geq 0} \bigcup_{s \in S} P_1^r(q, s), \quad \bigcup_{r \geq 0} \bigcup_{s \in S} P_2^r(q, s), \quad \bigcup_{r \geq 0} \bigcup_{s \in S} P_5^r(q, s), \quad \bigcup_{r \geq 0} \bigcup_{s \in S} P_6^r(q, s)$$

is bounded by $(tm)^{O(d^2 k)}$.

By Corollary 8.4 we are able to bound the number of all possible sets $\bigcup_{r \geq 0} \bigcup_{s \in S} P_7^r(q, s)$ as $(tm)^{O(d^2 k)}$. And because this bound is proven using Theorem 6.3, then it applies to the dual range space, and we also bound the number of possible sets in $\bigcup_{r \geq 0} \bigcup_{s \in S} P_8^r(q, s)$ as $(tm)^{O(d^2 k^2)}$. The k^2 term arises because we consider $\Theta(k^2)$ pairs q_i, q_t for predicate P_8 . So, ultimately,

$$2^t \leq (tm)^{O(d^2 k^2)} \implies t = O(d^2 k^2 \log dkm).$$

Similarly, we can bound the shattering dimension δ ,

$$t^\delta \leq (tm)^{O(d^2 k^2)} \implies \delta = O(d^2 k^2 \log m). \quad \square$$

9 Lower Bounds

Our lower bounds are constructed in the simplified setting that either $k = 1$ or $m = 1$, i.e., either the ground set or the curves defining the metric ball consist of one vertex only. In this case, all of our considered distance measures (except for one direction of the directed Hausdorff distance) are equal:

Lemma 9.1 *Let $p \in \mathbb{R}^d$, $q \in \mathbb{X}_k^d$, let $r = \max_{s \in V(q)} \|s - p\|$. Let $d_{dH}(p, q)$ be the Hausdorff distance between $V(p)$ and $V(q)$. It holds that*

$$r = d_{dH}(q, p) = d_{\overleftarrow{H}}(q, p) = d_{dF}(q, p) = d_F(q, p) = d_{wF}(q, p) = d_{dH}(p, q).$$

Proof In the discrete case we interpret $q \in \mathbb{X}_k^d$ as an ordered or unordered sequence of points in \mathbb{R}^d . In this case, the proof follows directly from definitions (Sect. 2). In the continuous case we interpret $q \in \mathbb{X}_k^d$ as a continuous polygonal curve. In this case, the proof follows directly from the definitions and from the convexity of the Euclidean ball of radius r centered at the point p . If and only if all vertices of q are contained in this ball, the distance is less or equal r .

Because of the above lemma, any lower bound that we prove for the Hausdorff distance in the discrete setting automatically extends to the other distance measures.

Lemma 9.2 *Let $\mathcal{R}_{dH,k}$ be the set of all balls, under the Hausdorff distance, centered at point sets in \mathbb{X}_k^2 . The VC dimension of the range space $(\mathbb{X}_m^2, \mathcal{R}_{dH,k})$ is $\Omega(k)$.*

Proof The intuition of our proof is as follows. We construct a set of k points in \mathbb{R}^2 that can be shattered by the ranges in $\mathcal{R}_{dH,k}$. The basic idea is that the ranges behave like convex polygons with k facets. In particular, the set of points contained inside the range centered at a curve q , is equal to the intersection of a set of equal-size Euclidean balls centered at the vertices of q .

Concretely, we place a set P of $k \geq 4$ points evenly spaced on a unit circle centered at the origin, see Fig. 6. Let $R > 2$ be a parameter of the construction. For representing any subset of P we construct q using k vertices (in any order) placed on the origin-centered circle of radius $R - 1$. In particular, we can force any $p_0 \in P$ to be excluded from the metric ball under the Hausdorff distance of a fixed radius

$$R' \in \left[\sqrt{R^2 - 2(R - 1) \left(1 - \cos \frac{2\pi}{k} \right)}, R \right),$$

by placing a vertex on the line through the origin that contains p_0 and by adding this vertex to the vertex set of q . Using the k vertices in q we can specifically exclude any subset of up to k points from P by such a construction, and by placing a vertex of q at the origin we will not exclude any points. Hence any set P on the unit circle of size k can be shattered. □

Lemma 9.3 *Let $\mathcal{R}_{dH,k}$ be the set of all balls, under the Hausdorff distance, centered at discrete point sets in \mathbb{X}_k^2 . The VC dimension of the range space $(\mathbb{X}_m^2, \mathcal{R}_{dH,k})$ is $\Omega(\log m)$.* □

Proof Lemma 9.2 and [30, Lem. 5.18], which bounds the VC dimension of the dual range space as a function of the VC dimension of the primal space, imply the theorem. □

Theorem 9.4 *The VC dimension of the range spaces $(\mathbb{X}_m^2, \mathcal{R}_{dF,k})$, $(\mathbb{X}_m^2, \mathcal{R}_{dH,k})$, $(\mathbb{X}_m^2, \mathcal{R}_{wF,k})$, $(\mathbb{X}_m^2, \mathcal{R}_{F,k})$, and $(\mathbb{X}_m^2, \mathcal{R}_{H,k})$ is $\Omega(\max(k, \log m))$.*

Proof It follows by applying Lemmas 9.2 and 9.3 together with Lemma 9.1. □

Lemma 9.5 *Let $\mathcal{R}_{dH,k}$ be the set of all balls, under the Hausdorff distance, centered at point sets in \mathbb{X}_k^d . For $d \geq 4$, the VC dimension of the range space $(\mathbb{X}_m^d, \mathcal{R}_{dH,k})$ is $\Omega(dk \log k)$.*

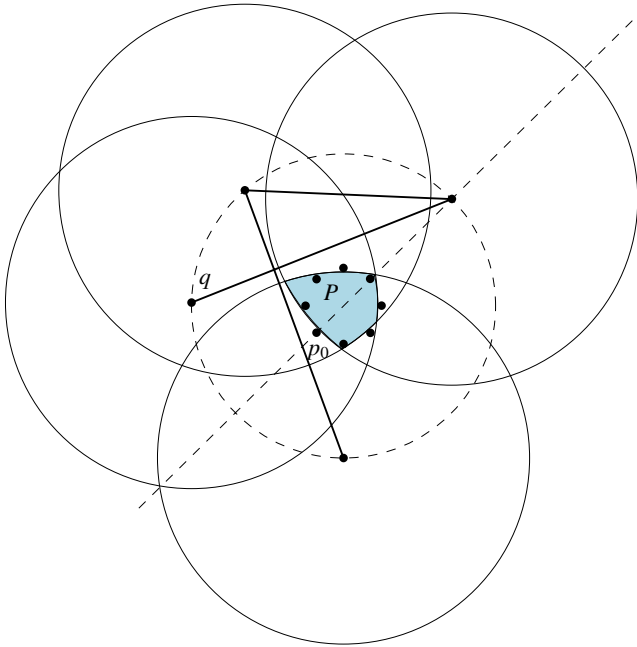


Fig. 6 A curve q with metric ball of radius R' containing a subset of P . The shaded area is the set of points that are contained inside the metric ball

Proof As in the proof of Lemma 9.2, our construction is set in the simplified setting where $m = 1$, i.e., the ground set corresponds to points in \mathbb{R}^d . We now show the theorem by reducing it to a recent lower bound of Csikós et al. [18] which is $\Omega(dk \log k)$ for a related range space for $d \geq 4$. This is defined on a ground set $P \subseteq \mathbb{R}^d$ with ranges \mathcal{R}_k defined so that each range $R \in \mathcal{R}_k$ is the intersection of k halfspaces. Recall that the construction in the proof of Lemma 9.2 used the fact that for $d = 2$ the ranges behave like convex polygons. We can observe a similar behavior in higher dimensions. In particular, Lemma 9.1 implies that any range in $\mathcal{R}_{dH,k}$ corresponds to the intersection of k balls in \mathbb{R}^d (centered at vertices of q). Observe that for a sufficiently large fixed radius R , for any point set $P \subseteq \mathbb{R}^d$, and for any halfspace H , we can find a ball of radius R which has the same inclusion properties as H . Finally, the lower bound of Csikós et al. [18] shows that there exists a set P of $\kappa = \Omega(dk \log k)$ points which can be shattered by such ranges. \square

Lemma 9.6 Let $\mathcal{R}_{dH,k}$ be the set of all balls, under the Hausdorff distance, centered at point sets in \mathbb{X}_k^d . For $d \geq 4$, the VC dimension of the range space $(\mathbb{X}_m^d, \mathcal{R}_{dH,k})$ is $\Omega(\log dm)$.

Proof Lemma 9.5 and [30, Lem. 5.18], which bounds the VC dimension of the dual range space as a function of the VC dimension of the primal space, imply the theorem. \square

Theorem 9.7 For $d \geq 4$, the VC dimension of $(\mathbb{X}_m^d, \mathcal{R}_{dF,k})$, $(\mathbb{X}_m^d, \mathcal{R}_{dH,k})$, $(\mathbb{X}_m^d, \mathcal{R}_{wF,k})$, $(\mathbb{X}_m^d, \mathcal{R}_{F,k})$, and $(\mathbb{X}_m^d, \mathcal{R}_{H,k})$ is

$$\Omega(\max(dk \log k, \log dm)).$$

Proof It follows by applying Lemmas 9.5 and 9.6 together with Lemma 9.1. \square

10 Implications

In this section we demonstrate that bounds on the VC dimension for the range space defined by metric balls on curves immediately imply various results about prediction and statistical generalization over the space of curves. In the following consider a range space (X, \mathcal{R}) with a ground set X of curves, where \mathcal{R} are the ranges corresponding to metric balls for some distance measure we consider, and the VC dimension is bounded by ν .

This section discusses accuracy bounds that depend directly on the size $n = |X|$ and the VC dimension ν . We will assume that X is a random sample of some much larger set X_{big} or an unknown continuous generating distribution μ . Under the randomness in this assumed sampling procedure, there is a probability of failure δ that often shows up in these bounds, but is minor since it shows up as $\log(1/\delta)$. The following discussion leverages the concepts of ε -nets and ε -samples. The former (ε -nets) are samples which satisfy the property that if a range is heavy (contains an ε -fraction of the data) then the sample contains at least one point in that range; a sample of size $O((\nu/\varepsilon) \log(\nu/\varepsilon\delta))$ is sufficient [33]. The latter (ε -samples) are samples which satisfy that each range's density is approximated within an additive ε -error; a sample of size $O((\nu - \log \delta)/\varepsilon^2)$ is sufficient [39].

These bounds often take two closely-linked forms. First, given a limited set X from an unknown μ , then how accurate is a query or a prediction made using only X . Second, given the ability to draw samples (at a cost) from an unknown distribution μ , how many are required so that the prediction on the set of samples X has bounded prediction error. Upper bounds on ν imply pessimistic bounds on the accuracy or the required size for a sample.

Such large data sets of curves are now commonplace in many structured data applications. For instance, the millions of ride-sharing trips taken every day, or the GPS traces Apple and Google and others collect on users' phones, or the tracking of migrating animals. Because this data has a complex structure, and each associated curve may be large (i.e., m is large), it is not clear how well analyses on families of such curves can provably generalize to predict new data. The theme of the following results, as implied by our above VC dimension results, is that if these families of curves are only inspected with or queried with curves with a small number of segments (i.e., k is small), then the VC dimension of the associated range space $\nu = O(k \log km)$ or $O(k^2 \log km)$ is small, and that such analyses generalize well. We show this in several concrete examples.

Approximate range counting on curves. Given a large set of curves X (of potentially very large complexity m) and a query curve q (with smaller complexity k) we would

like to approximate the number of curves nearby q . For instance, we restrict X to historical queries at a certain time of day and query with the planned route q , and would like to know the chance of finding a carpool. VC dimension ν of the metric balls shows up directly in two analyses. First, if we assume that X has been chosen from an unknown distribution, i.e., $X \sim \mu$ where μ is a much larger unknown distribution (but the real one), then we can estimate the accuracy of the fraction of all curves in this range within additive error $O(\sqrt{(1/|X|)(\nu + \log(1/\delta))})$. On the other hand, if we assume that X is a fixed input set which is too large to conveniently query, we can sample a subset $S \subset X$ of size $O((1/\varepsilon^2)(\nu + \log(1/\delta)))$ and know that the estimate for the fraction of curves from S in that range is within additive ε error of the fraction from X . Such sampling techniques have a long history in traditional databases [40], and have more recently become important when providing online estimates during a long query processing time as incrementally increasing size subsets are considered [3]. Ours provides the first formal analysis of these results for queries over curves. Moreover, the finite bound on VC dimension of these problems also implies [17] that there is a linear size data structure which can answer exact range queries in sublinear time.

Density estimation of curves. A related task in generalization to new curves is density estimation. Consider a large set of curves X which represent a larger unknown distribution μ that models a distribution of curves; we want to understand how unusual a new curve q would be, given we have not yet seen exactly the same curve before. One option is to use the distance to the (k th) nearest neighbor curve in X , or a bit more robust option is to choose a radius r and count how many curves are within that radius (e.g., the approximate range counting results above).

Alternatively, for $X \subset \mathbb{M}$, consider now a kernel density estimate $\text{KDE}_X: \mathbb{M} \rightarrow \mathbb{R}$ defined by

$$\text{KDE}_X(p) = \frac{1}{n} \sum_{p \in P} K(x, p)$$

with kernel $K(x, p) = \exp(-d(x, p)^2)$ (where d is some distance of choice among curves, e.g., d_F). The kernel is defined so that each superlevel set $K_x^\tau = \{p \in \mathbb{M} \mid K(x, p) \geq \tau\}$ corresponds to some range $R \in \mathcal{R}$ such that $R \cap X = K_x^\tau \cap X$. Then a random sample $S \subset X$ of size $O((1/\varepsilon^2)(\nu + \log(1/\delta)))$ satisfies $\|\text{KDE}_X - \text{KDE}_S\|_\infty \leq \varepsilon$ [35]. Thus, again the VC dimension ν of the metric balls directly influences these estimates accuracy, and for query curves with small complexity k the bound is quite reasonable.

Sample complexity for classification of curves. Now consider the problem of classifying curves representing trajectories of people or animals. For instance, with individuals who enable GPS on their cell phone they can label some work-to-home trajectories (as $\chi(x) = +1$) or as other trips ($\chi(x) = -1$). Then on unlabeled trips we can potentially predict which are work-to-home trajectories to build traffic and commute time models without manually labeling all routes. Similar tasks may be useful for normal ($\chi(x) = +1$) versus nefarious ($\chi(x) = -1$) activities when tracking people in an airport or a hostile zone. In each of these cases we may either have a very large number

of labeled instances, and may want to sample them to some manageable size, or we may only have a limited number of samples, and want to know the accuracy to trust based on the sample size. All of these bounds are controlled by the VC dimension of the family of classifiers used to make the prediction. For trajectories, a sensible family of classifiers would be the ranges \mathcal{R} defined by metric balls.

That is, consider some labeling function $\chi: X \rightarrow \{-1, +1\}$; now we say a range $R \subset \mathcal{R}$ misclassifies an object $x \in X$ if $x \in R$ and $\chi(x) = -1$ or $x \notin R$ and $\chi(x) = +1$. If there exists a range $R \subset \mathcal{R}$ such that all $x \in X \cap R$ have $\chi(x) = +1$ and all $x' \in X \setminus R$ have $\chi(x') = -1$, we say such a range *perfectly separates* (X, χ) . Then a random sample $Y \subset X$ of size $O((v/\varepsilon) \log(v/\varepsilon\delta))$ [33] ensures that, with probability at least $1 - \delta$, any range $R' \subset \mathcal{R}$ which perfectly separates (Y, χ) misclassifies at most εn points in X .

Consider a random sample $Y \subset X$ of size $O((1/\varepsilon^2)(v + \log(1/\delta)))$. For any range $R \subset \mathcal{R}$, if the fraction of points in Y is $|R \cap Y|/|Y| = \eta$, then with probability at least $1 - \delta$, the fraction of points in X is $|R \cap X|/|X| \in [\eta - \varepsilon, \eta + \varepsilon]$; that is, it is off by at most an ε -fraction [31,39]. If there is a labeling $\chi: X \rightarrow \{-1, +1\}$, this notably includes the range $R \in \mathcal{R}$ which misclassifies the least points (there may not be a perfect separator). Hence a random sampling ensures at most an ε -fraction more misclassified points are included in an estimate derived from this sample. Indeed, the RBF kernel $K(x, p) = \exp(-d(x, p)^2)$ defined above implies standard mechanism like kernel SVM or kernel perceptron [42] can be used to build classifiers, and together these bounds induce misclassification [39] and margin approximation bounds [35]. The small VC dimension v implies they will generalize well.

Funding Open Access funding enabled and organized by Projekt DEAL. Anne Driemel received funding from the German Research Foundation (DFG) and the Netherlands Organization for Scientific Research (NWO). André Nusser received funding from the Max Planck Society for the Advancement of Science. Jeff Phillips received funding from the National Science Foundation (NSF) and Visa Research; part of the work was completed while visiting the Simons Institute for Theory of Computing. Ioannis Psarros received funding from the State Scholarships Foundation (IKY); this research is co-financed by Greece and the European Union (European Social Fund—ESF).

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

References

1. Afshani, P., Driemel, A.: On the complexity of range searching among curves (2017). [arXiv:1707.04789](https://arxiv.org/abs/1707.04789)
2. Afshani, P., Driemel, A.: On the complexity of range searching among curves. In: 29th Annual ACM-SIAM Symposium on Discrete Algorithms (New Orleans 2018), pp. 898–917. SIAM, Philadelphia (2018)

3. Agarwal, S., Mozafari, B., Panda, A., Milner, H., Madden, S., Stoica, I.: BlinkDB: queries with bounded errors and bounded response times on very large data. In: 8th ACM European Conference on Computer Systems (Prague 2013). ACM, New York (2013)
4. Akama, Y., Irie, K., Kawamura, A., Uwano, Y.: VC dimensions of principal component analysis. *Discret. Comput. Geom.* **44**(3), 589–598 (2010)
5. Alt, H., Behrends, B., Blömer, J.: Approximate matching of polygonal shapes. *Ann. Math. Artif. Intell.* **13**(3–4), 251–265 (1995)
6. Alt, H., Godau, M.: Computing the Fréchet distance between two polygonal curves. *Int. J. Comput. Geom. Appl.* **5**(1–2), 75–91 (1995)
7. Anthony, M., Bartlett, P.L.: *Neural Network Learning: Theoretical Foundations*. Cambridge University Press, Cambridge (1999)
8. Astefanoaei, M., Cesaretti, P., Katsikouli, P., Goswami, M., Sarkar, R.: Multi-resolution sketches and locality sensitive hashing for fast trajectory processing. In: 26th ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems (Seattle 2018), pp. 279–288. ACM, New York (2018)
9. Baldus, J., Bringmann, K.: A fast implementation of near neighbors queries for Fréchet distance (GIS Cup). In: 25th ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems (Redondo Beach 2017), # 99. ACM, New York (2017)
10. de Berg, M., Cook IV, A.F., Gudmundsson, J.: Fast Fréchet queries. *Comput. Geom.* **46**(6), 747–755 (2013)
11. de Berg, M., Mehrabi, A.D.: Straight-path queries in trajectory data. In: Algorithms and Computation—9th International Workshop (Dhaka 2015). Lecture Notes in Computer Science, vol. 8973, pp. 101–112. Springer, Cham (2015)
12. Blumer, A., Ehrenfeucht, A., Haussler, D., Warmuth, M.K.: Learnability and the Vapnik–Chervonenkis dimension. *J. Assoc. Comput. Mach.* **36**(4), 929–965 (1989)
13. Bringmann, K., Künnemann, M., Nusser, A.: Walking the dog fast in practice: algorithm engineering of the Fréchet distance. In: 35th International Symposium on Computational Geometry (Portland 2019). Leibniz Int. Proc. Inform., vol. 129, # 17. Leibniz-Zent. Inform., Wadern (2019)
14. Brönnimann, H., Goodrich, M.T.: Almost optimal set covers in finite VC-dimension. *Discret. Comput. Geom.* **14**(4), 463–479 (1995)
15. Buchin, K., Diez, Y., van Diggelen, T., Meulemans, W.: Efficient trajectory queries under the Fréchet distance (GIS Cup). In: 25th ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems (Redondo Beach 2017), # 101. ACM, New York (2017)
16. Ceccarello, M., Driemel, A., Silvestri, F.: FRESH: Fréchet similarity with hashing. In: Algorithms and Data Structures—16th International Symposium (Edmonton 2019). Lecture Notes in Computer Science, vol. 11646, pp. 254–268. Springer, Cham (2019)
17. Chazelle, B., Welzl, E.: Quasi-optimal range searching in spaces of finite VC-dimension. *Discret. Comput. Geom.* **4**(5), 467–489 (1989)
18. Csikós, M., Kupavskii, A., Mustafa, N.H.: Optimal bounds on the VC-dimension (2018). [arXiv:1807.07924](https://arxiv.org/abs/1807.07924)
19. Csikós, M., Mustafa, N.H., Kupavskii, A.: Tight lower bounds on the VC-dimension of geometric set systems. *J. Mach. Learn. Res.* **20**, # 81 (2019)
20. Driemel, A., Krivošija, A., Sohler, Ch.: Clustering time series under the Fréchet distance. In: 27th Annual ACM-SIAM Symposium on Discrete Algorithms (Arlington 2016), pp. 766–785. ACM, New York (2016)
21. Driemel, A., Silvestri, F.: Locally-sensitive hashing of curves. In: 33rd International Symposium on Computational Geometry (Brisbane 2017). Leibniz Int. Proc. Inform., vol. 77, # 37. Leibniz-Zent. Inform., Wadern (2017)
22. Dütsch, F., Vahrenhold, J.: A filter-and-refinement-algorithm for range queries based on the Fréchet distance (GIS Cup). In: 25th ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems (Redondo Beach 2017), # 100. ACM, New York (2017)
23. Eiter, T., Mannila, H.: Computing discrete Fréchet distance. Tech. Rep. CD-TR 94/64. Technische Universität, Wien (1994)
24. Emiris, I.Z., Psarros, I.: Products of Euclidean metrics and applications to proximity questions among curves. In: 34th International Symposium on Computational Geometry (Budapest 2018). Leibniz Int. Proc. Inform., vol. 99, # 37. Leibniz-Zent. Inform., Wadern (2018)

25. Fréchet, M.M.: Sur quelques points du calcul fonctionnel. *Rendiconti del Circolo Matematico di Palermo* **22**, 1–72 (1906)
26. Gilbers, A., Klein, R.: A new upper bound for the VC-dimension of visibility regions. *Comput. Geom.* **47**(1), 61–74 (2014)
27. Goldberg, P.W., Jerrum, M.R.: Bounding the Vapnik–Chervonenkis dimension of concept classes parameterized by real numbers. *Mach. Learn.* **18**(2–3), 131–148 (1995)
28. Gudmundsson, J., Horton, M.: Spatio-temporal analysis of team sports. *ACM Comput. Surv.* **50**(2), # 22 (2017)
29. Gudmundsson, J., Smid, M.: Fast algorithms for approximate Fréchet matching queries in geometric trees. *Comput. Geom.* **48**(6), 479–494 (2015)
30. Har-Peled, S.: *Geometric Approximation Algorithms*. Mathematical Surveys and Monographs, vol. 173. American Mathematical Society, Providence (2011)
31. Har-Peled, S., Sharir, M.: Relative (p, ϵ) -approximations in geometry. *Discret. Comput. Geom.* **45**(3), 462–496 (2011)
32. Hausdorff, F.: *Grundzüge der Mengenlehre*. Veit, Leipzig (1914)
33. Haussler, D., Welzl, E.: ϵ -nets and simplex range queries. *Discret. Comput. Geom.* **2**(2), 127–151 (1987)
34. Huang, L., Jiang, S.H.-C., Li, J., Wu, X.: ϵ -coresets for clustering (with outliers) in doubling metrics. In: 59th Annual IEEE Symposium on Foundations of Computer Science (Paris 2018), pp. 814–825. IEEE Computer Soc., Los Alamitos (2018)
35. Joshi, S., Kommaraju, R.V., Phillips, J.M., Venkatasubramanian, S.: Comparing distributions and shapes using the kernel distance. In: 27th Annual Symposium on Computational Geometry (Paris 2011), pp. 47–56. ACM, New York (2011)
36. Karpinski, M., Macintyre, A.: Polynomial bounds for VC dimension of sigmoidal and general Pfaffian neural networks. *J. Comput. Syst. Sci.* **54**(1), 169–176 (1997)
37. Langetepe, E., Lehmann, S.: Exact VC-dimension for L_1 -visibility of points in simple polygons (2017). [arXiv:1705.01723](https://arxiv.org/abs/1705.01723)
38. Laurila, J.K., Gatica-Perez, D., Aad, I., Blom, J., Bornet, O., Do, T.-M.-T., Dousse, O., Eberle, J., Miettinen, M.: From big smartphone data to worldwide research: the Mobile Data Challenge. *Pervas. Mob. Comput.* **9**(6), 752–771 (2013)
39. Li, Y., Long, P.M., Srinivasan, A.: Improved bounds on the sample complexity of learning. *J. Comput. Syst. Sci.* **62**(3), 516–527 (2001)
40. Olken, F.: *Random Sampling from Databases*. PhD thesis, University of California at Berkeley (1993). <https://dsf.berkeley.edu/papers/UCB-PhD-olken.pdf>
41. Sauer, N.: On the density of families of sets. *J. Comb. Theory Ser. A* **13**, 145–147 (1972)
42. Schölkopf, B., Smola, A.J.: *Learning with Kernels*. MIT Press, Cambridge (2002)
43. Shelah, S.: A combinatorial problem; stability and order for models and theories in infinitary languages. *Pac. J. Math.* **41**, 247–261 (1972)
44. Valtr, P.: Guarding galleries where no point sees a small area. *Israel J. Math.* **104**, 1–16 (1998)
45. Vapnik, V.N.: *Statistical Learning Theory*. Wiley, New York (1998)
46. Vapnik, V.N., Chervonenkis, A.Ya.: On the uniform convergence of relative frequencies of events to their probabilities. *Theor. Probab. Appl.* **16**(2), 264–280 (1971). (in Russian)
47. Werner, M., Oliver, D.: ACM SIGSPATIAL GIS Cup 2017: range queries under Fréchet distance. *SIGSPATIAL Spec.* **10**(1), 24–27 (2018)
48. Zheng, F., Kaiser, T.: *Digital Signal Processing for RFID*. Wiley Series on Information and Communication Technology. Wiley, Chichester (2016)

Authors and Affiliations

Anne Driemel¹ · André Nusser² · Jeff M. Phillips³ · Ioannis Psarros^{1,4} 

Anne Driemel
driemel@cs.uni-bonn.de

André Nusser
anusser@mpi-inf.mpg.de

Jeff M. Phillips
jeffp@cs.utah.edu

Ioannis Psarros
ipsarros@di.uoa.gr ; ipsarros@cs.uni-bonn.de

- ¹ Present Address: University of Bonn, Endenicher Allee 19a, 53115 Bonn, Germany
- ² Max Planck Institute for Informatics & Graduate School of Computer Science, Saarland Informatics Campus, Campus E1 4, 66123 Saarbrücken, Germany
- ³ University of Utah, 50 S. Central Campus Dr., Salt Lake City, UT 84112, USA
- ⁴ National and Kapodistrian University of Athens, 15784 Panepistimiopolis, Greece