

A Range Space with Constant VC Dimension for All-pairs Shortest Paths in Graphs

Alane Marie de Lima¹ Murilo V. G. da Silva² André L. Vignatti²

¹Federal University of Technology – Paraná, Brazil

²Federal University of Paraná, Brazil

Submitted: December 2022	Reviewed: April 2023	Revised: June 2023
Reviewed: July 2023	Revised: August 2023	Accepted: August 2023
Final: September 2023	Published: September 2023	
Article type: Regular Paper	Communicated by: G. Liotta	

Abstract. Let G be an undirected graph with non-negative edge weights and let S be a subset of its shortest paths such that, for every pair (u, v) of distinct vertices, S contains exactly one shortest path between u and v . In this paper we define a range space associated with S and prove that its VC dimension is 2. As a consequence, we show a bound for the number of shortest paths trees required to be sampled in order to solve a relaxed version of the All-pairs Shortest Paths problem (APSP) in G . In this version of the problem we are interested in computing all shortest paths with a certain “importance” at least ε . Given any $0 < \varepsilon, \delta < 1$, we propose a $\mathcal{O}(m + n \log n + (\text{diam}_{V(G)})^2)$ sampling algorithm that outputs with probability $1 - \delta$ the (exact) distance and the shortest path between every pair of vertices (u, v) that appears as subpath of at least a proportion ε of all shortest paths in the set S , where $\text{diam}_{V(G)}$ is the vertex-diameter of G . The bound that we obtain for the sample size depends only on ε and δ , and do not depend on the size of the graph.

1 Introduction

The All-pairs Shortest Path (APSP) is the problem of computing a path with the minimum length between every pair of vertices in a weighted graph. The APSP problem is very well studied and there has been recent results for a variety of assumptions for the input graph (directed/undirected,

Research supported by CAPES (Proc. Proc. 88882.461738/2019-01) and CNPq (Proc. 428941/2016-8 and 420079/2021-1).

E-mail addresses: alanelima@utfpr.edu.br (Alane Marie de Lima) murilo@inf.ufpr.br (Murilo V. G. da Silva) vignatti@inf.ufpr.br (André L. Vignatti)



This work is licensed under the terms of the [CC-BY](https://creativecommons.org/licenses/by/4.0/) license.

integer/real edge weights, etc) [27, 5, 9, 4]. In this paper we assume that the input is an undirected graph G with n vertices and m edges with non-negative weights.

In our scenario, the fastest known exact algorithms are the algorithm proposed by Williams (2014) [27], which runs in $\mathcal{O}\left(\frac{n^3}{2^{c\sqrt{\log n}}}\right)$ time, for some constant $c > 0$, and by Pettie and Ramachandram (2002) [18] for the case of sparse graphs, which runs in $\mathcal{O}(nm \log \alpha(m, n))$ time, where $\alpha(m, n)$ is the Tarjan’s inverse-Ackermann function. If no assumption is taken about the sparsity of the graph, then it is an open question whether the APSP problem can be solved in *strictly* subcubic time, i.e. $\mathcal{O}(n^{3-c})$, for any $c > 0$, even when the edge weights are natural numbers.

Recent results in fine-grained complexity indicate that the complexity time for the APSP is tight [22, 1, 2], reinforcing the hypothesis that there is no strictly subcubic algorithm for such task [25]. Since the exact computation of this version is expensive for large graphs, especially the dense ones, it is natural dealing with alternative versions of the problem, whether they are approximate [8, 21] or applied to restricted scenarios [24]. In this paper, we follow this line of work, dealing with a relaxation of the problem in the sense that the classical APSP is a special case for a given adjustable parameter. More specifically, we aim to compute, with high probability, all the shortest paths that meet a certain “importance” requirement. The idea is that the probability of sampling a shortest path P is higher when a large number of shortest paths from some set of *canonical shortest paths* has P as a subpath. The precise definition of this measure is given in Section 2.

Let S be a subset of the shortest paths in G that contains exactly one shortest path between a pair of distinct vertices, for all $(u, v) \in V^2$. In this relaxed version of the APSP, given constant parameters $0 < \varepsilon, \delta < 1$, we propose a sampling algorithm that outputs, with probability at least $1 - \delta$, the (exact) distance and a shortest path between every pair of vertices that admits a shortest path that appears as subpath of at least a proportion ε of all shortest paths in the set S . The central idea of the algorithm is to sample roots of shortest paths trees. In order to give a bound for the sample size that is sufficient to meet the input parameters, we use sample complexity tools, namely, Vapnik–Chervonenkis (VC) dimension theory and the ε -net theorem. We define a range space associated with S in the graph G . One of the main results that we prove is that the VC dimension of such range space is 2 and that the bound for the sample size is $r = \lceil \frac{c}{\varepsilon} (2 \ln(\frac{1}{\varepsilon}) + \ln \frac{1}{\delta}) \rceil$, where c is a constant around $\frac{1}{2}$ [15]. This result is interesting, since it does not depend neither on the size of the input n , which is the case if one uses standard union-bound techniques, nor on the topological structure of the graph that may vary with n in many cases. As a consequence of this bound for the sample size, we obtain a sampling algorithm for our problem with running time $\mathcal{O}(m + n \log n + (\text{Diam}_V(G))^2)$, where $\text{Diam}_V(G)$ is the vertex-diameter of the input graph (i.e. the maximum number of vertices in a shortest path in G), for any constant ε .

If one sets ε as a function of n , in the limit case, when $\varepsilon(n) = \frac{1}{n(n-1)}$, our algorithm solves – with high probability – the classical APSP problem, but with time complexity exceeding the running time of the exact algorithms from the literature [28, 18]. However, it is still an interesting problem to know for which functions $\varepsilon(n)$ we still have a strictly subcubic sampling algorithm. We show that our algorithm runs in $\mathcal{O}(n^{3-c})$ time if $\varepsilon(n)$ is any $\Omega\left(\frac{W_0(n')}{n'}\right)$ function, where $n' = n^{1-c}$ (for a constant $c > 0$) and $W_0(n')$ is the branch 0 of the Lambert-W function defined for $n' \geq 0$, a non-algebraic value such that $W_0(n') = \ln n' - \ln \ln n' + \Theta\left(\frac{\ln \ln n'}{\ln n'}\right)$, which holds for $n' \geq e$.

2 Shortest Paths, Canonical Paths, and Shortest Paths Trees

Let $G = (V, E)$ be an undirected graph, with $n = |V|$ and $m = |E|$, and let ω be a function of edge *weights* from E to an enumerable subset of $\mathbb{R}_{\geq 0}$. W.l.o.g., we assume that G is connected, since our results can be applied to the connected components when a graph is disconnected. Even though G is undirected, for convenience we use the notation (u, v) for an edge of G . A *path* is a sequence of vertices $P = (v_1, v_2, \dots, v_k)$ such that $v_i \neq v_{i+1}$ and $(v_i, v_{i+1}) \in E$, for $1 \leq i < k$. If $u = v_1$ and $v = v_k$, such path is referred to as a (u, v) -*path*. We define E_P as the set of edges of P . The *shortest path* from u to v in G is the (u, v) -path such that the sum of the weights of the edges in E_P is minimized. In this case we denote such value $d(u, v)$, also called the *distance* from u to v .

The set of all shortest paths from u to v in G is denoted \mathcal{C}_{uv} . For a given path $P \in \mathcal{C}_{uv}$, let $\text{Inn}(P)$ be the set of *inner* vertices of P , that is, $\text{Inn}(P) = \{w \in P : w \notin \{u, v\}\}$. Consider a shortest (u, v) -path P , and let u' and v' be two vertices of P , with u' closer to u and v' closer to v . The subpath of P starting in u' and ending in v' is called a (u', v') -*subpath* of P . The (immediate) *predecessor* of v in a shortest (u, v) -path P , denoted $\text{pred}_P(v)$, is the vertex $w \in \text{Inn}(P)$ such that $(w, v) \in E_P$. The *diameter* of G , denoted Diam_G , is the size of the largest shortest path in G . The *vertex-diameter*, denoted $\text{Diam}_V(G)$, is the maximum number of vertices in a shortest path of G .

Let $\sigma : V \rightarrow \{1, \dots, n\}$ be an arbitrary vertex ordering of G . Consider the set of shortest paths $\mathcal{L}_{uv} = \{P \in \mathcal{C}_{uv} : \sigma(\text{pred}_P(v)) \text{ is minimum}\}$. Note that there is only one vertex w that satisfies the property “ $\sigma(\text{pred}_P(v))$ is minimum”, so even if there are several paths in \mathcal{L}_{uv} , the last edge (w, v) is the same for all of them. Next, we introduce the definition of a *canonical path* with respect to σ .

Definition 2.1 (Canonical path (CP)) Consider a pair of vertices $(u, v) \in V^2$ in G . The canonical path (CP) from u to v , denoted P , is recursively defined as the shortest path in \mathcal{C}_{uv} such that

case 1: $|\mathcal{L}_{uv}| = 1$. Then $P \in \mathcal{L}_{uv}$ is the canonical path from u to v .

case 2: $|\mathcal{L}_{uv}| > 1$. Let w be the (unique) predecessor of v in the shortest paths of \mathcal{L}_{uv} . Then, the canonical path from u to v corresponds to the canonical path from u to w plus the edge (w, v) .

Fact 1 Given a pair of vertices $(u, v) \in V^2$, the CP from u to v exists and it is unique.

To see that Fact 1 holds, note that at each recursive step, there is only one vertex w satisfying the property that defines \mathcal{L}_{uv} , and there is only one canonical path from u to w . Besides, the recursion presented above always stop in the base case, since the distance between a pair of vertices in a recursive step is smaller than the distance of a pair of vertices analyzed in the previous step. The base is the one where there is only one (u, u') -subpath which is the shortest path from u to u' , for $u' \in \text{Inn}(P)$. Another important observation about canonical paths is that the canonical path from u to v is not necessarily the same as the canonical path from v to u .

A *shortest paths tree (SPT)* of a vertex u is a spanning tree of G such that the path from u to every other vertex of this tree is a shortest path in G . There might be many SPTs for a given vertex. In this paper we are interested in fixing one canonical SPT T_u , for every vertex u of G . More precisely, for a given (arbitrary) vertex ordering σ , the canonical SPT T_u is defined such that, for every vertex v , the shortest path from u to v in T_u is a canonical path. In Section 4.1 we give more details on the computation of T_u , but, briefly speaking, this tree is the one computed by a modification on Dijkstra’s algorithm where σ is used as a tie-breaking criterion. We also call T_u the *Dijkstra tree* of u .

A shortest path that starts at the root of a Dijkstra tree is also called a *branch* of G . More formally, given T_u , for every $v \neq u$, the shortest path from u to v is a branch, denoted \mathcal{B}_{uv} . In addition, every subpath of \mathcal{B}_{uv} is also a shortest path in G , and we denote such set of subpaths (including \mathcal{B}_{uv}) as $S(\mathcal{B}_{uv})$.

We denote $c(u, v)$ as the proportion of canonical shortest paths containing a shortest path between u and v as subpath. In order to formally define $c(u, v)$ we first need the following. Let t_{uv} be the number of canonical paths that contain a shortest path from u to v as subpath, defined as

$$t_{uv} = \sum_{(a,b) \in V^2, a \neq b} \mathbb{1}_{uv}(\mathcal{B}_{ab}),$$

where $\mathbb{1}_{uv}(\mathcal{B}_{ab})$ is the indicator function that returns 1 if there is some shortest path from u to v as subpath of the branch \mathcal{B}_{ab} (and 0 otherwise).

Definition 2.2 *Given a pair $(u, v) \in V^2$,*

$$c(u, v) = \frac{t_{uv}}{n(n-1)}, \quad \text{where } n = |V|.$$

2.1 Key Results on Canonical Paths

Before we present the main results of this paper in Section 3.1, we need first a key technical result concerning canonical paths. We show in Theorem 1 that any subpath of a canonical path is also a canonical path.

Lemma 1 *Given a pair of vertices $(u, v) \in V^2$, let P be the CP from u to v in G . If $|\mathcal{L}_{uv}| = 1$, then every subpath of P is also a CP.*

Proof: Let P' be a (u', v') -subpath of P . Suppose by contradiction that P' is not a CP. Let $Q' \neq P'$ be the shortest path $Q' = (u', \dots, v')$ in G which is the CP from u' to v' .

Case 1: $v' \neq v$. Let S_1 be a (u, u') -subpath and S_2 be a (v', v) -subpath, both from P . Let Q be the concatenation of S_1 , Q' , and S_2 . Note that P' and Q' have the same length (since both are shortest paths), and so does P and Q . Since P and Q have the same vertices from v' to v , then the predecessor of v in both paths is the same. Hence, P and Q are in \mathcal{L}_{uv} . But then $|\mathcal{L}_{uv}| > 1$, a contradiction.

Case 2: $v' = v$. Let w and w' be the predecessors of v in P' and Q' , respectively. Note that $w \neq w'$. Thus, since $\{w', v\}$ is the last edge of Q' , by the definition of CP, $\sigma(w') < \sigma(w)$. But then in the edge (w, v) of P , vertex w does not have the minimum index among all possible predecessors of v , contradicting the fact that P is a CP. \square

Lemma 2 *Given a pair of vertices $(u, v) \in V^2$, let P be the CP from u to v in G . Let w be the predecessor of v in P . Then the (u, w) -subpath of P is the CP from u to w .*

Proof: Let P' be the (u, w) -subpath of P . In the case of $|\mathcal{L}_{uv}| = 1$, then from Lemma 1 we have that P' is the CP from u to w . Otherwise, by Definition 2.1 (case 2) applied to P , it must hold that P' is the CP from u to w . \square

Lemma 3 *Given a pair of vertices $(u, v) \in V^2$, let P be the CP from u to v in G . Then for each $z \in \text{Inn}(P)$, the (u, z) -subpath of P is the CP from u to z .*

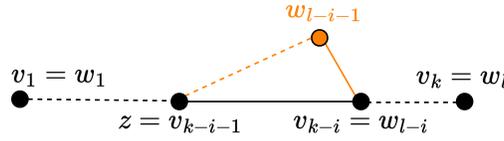


Figure 1: Illustration of vertex v_{k-i} in the shortest paths P (depicted in black color) and Z (depicted in orange color) in the proof of Lemma 4.

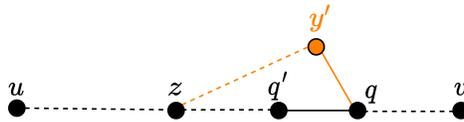


Figure 2: Illustration of the shortest path from z to q (in orange), denoted Q' , in the proof of Lemma 4.

Proof: Let P' be the (u, z) -subpath of P and w be the predecessor of v in P . We prove our claim by induction on the number of edges from z to v . The base case is the one where $z = w$ (i.e. P' is the (u, w) -subpath of P). This holds from Lemma 2.

Let z' be the predecessor of z in P and let P'' be the (u, z') -subpath of P . For the induction step, we show that if P' is CP from u to z , then P'' is the CP from u to z' .

By Definition 2.1 applied to P' , there are two cases to consider: $|\mathcal{L}_{uz}| = 1$ (case 1) and $|\mathcal{L}_{uz}| > 1$ (case 2). In case 1, by Lemma 1 applied to P' , the shortest path P'' must be the CP from u to z' . In case 2, by Definition 2.1 (case 2) applied to P' , the CP from u to z' is P'' . \square

Lemma 4 *Given a pair of vertices $(u, v) \in V^2$, let P be the CP from u to v in G . Then for each $z \in \text{Inn}(P)$, the (z, v) -subpath of P is the CP from z to v .*

Proof: Let Q be the (z, v) -subpath of P . We prove by contradiction supposing that Q is not the CP from z to v in G . Then there is a shortest path Y which is the CP from z to v in G . Consider the subpath of P from u to z concatenated with Y , and denote such concatenation as Z . Note that, even though the number of vertices of Q and Y may be different, the length of Q and Y is the same, since both are shortest paths. The same applies to P and Z .

Denote the vertices in P and Z as $P = (u = v_1, \dots, v = v_k)$ and $Z = (u = w_1, \dots, v = w_l)$. Let v_{k-i} be the vertex of P such that i is maximum, $0 \leq i < k$, and such that the following holds: for all $0 \leq j \leq i$, the vertex v_{k-j} in P is the same as the vertex w_{l-j} in Z (Figure 1). For simplicity, denote v_{k-i} as q , v_{k-i-1} as q' , and w_{l-i-1} as y' . Note that the edges in the (q, v) -subpaths of P and Z are the same, but (q', q) and (y', q) is not the same edge.

Let Q' and Y' be the (z, q) -subpaths of P and Y , respectively (Figure 2). Since we are assuming that Y is the CP from z to v in G , then by Lemma 3, Y' is the CP from z to q in G . Note that $Q' \neq Y'$ (since $Q \neq Y$), and hence, Q' is not the CP from z to q in G . Thus, $\sigma(q') > \sigma(y')$.

From Lemma 3 applied to P , the (u, q) -subpath of P is a CP. But this path is a shortest path such that q' is not the vertex with minimum index among all possible predecessors of q (recall that $\sigma(q') > \sigma(y')$), a contradiction. \square

Theorem 1 *Given a pair of vertices $(u, v) \in V^2$, let P be the CP from u to v in G . Then for each $(u', v') \in V^2$, the (u', v') -subpath of P is the CP from u' to v' in G .*

Proof: Let P' be the (u', v') -subpath of P . From Lemma 4, the (u', v) -subpath of P , denoted Q , is a CP. From Lemma 3, since Q is the CP from u' to v in G , then P' is the CP from u' to v' in G . □

3 Sample Complexity and VC Dimension

In sampling algorithms, typically the aim is the estimation of a certain quantity according to given parameters of quality and confidence using a random sample of size as small as possible. A central concept in sample complexity theory is the Vapnik–Chervonenkis Theory (VC dimension), in particular, the idea of finding an upper bound for the VC dimension of a class of binary functions related to the sampling problem at hand. In our context, for instance, we may consider a binary function that takes a branch and outputs 1 if such branch contains a shortest path for a given set. Generally speaking, from the upper bound for the VC dimension of the given class of binary functions we can derive an upper bound to the sample size for the sampling algorithm.

We present in this section the main definitions and results from sample complexity theory used in this paper. An in-depth exposition of the VC dimension theory and the ε -net theorem can be found in the books of Shalev-Schwartz and Ben-David (2014) [23], Mitzenmacher and Upfal (2017) [16], Anthony and Bartlett (2009) [3], and Mohri et al. (2012) [17].

Definition 3.1 (Range Space) *A range space is a pair $\mathcal{R} = (U, \mathcal{I})$, where U is a domain (finite or infinite) and \mathcal{I} is a collection of subsets of U , called ranges.*

For a given $S \subseteq U$, the *projection* of \mathcal{I} on S is the set $\mathcal{I}_S = \{S \cap I : I \in \mathcal{I}\}$. If $|\mathcal{I}_S| = 2^{|S|}$ then we say S is *shattered* by \mathcal{I} . Consider, for example, $S = \{1, 2, 3\}$ and the ranges $I_1 = \{1, 2, 5\}$ and $I_2 = \{1, 3, 4\}$. Then we have $S \cap I_1 = \{1, 2\}$ and $S \cap I_2 = \{1, 3\}$.

The VC dimension of a range space is the size of the largest subset S that can be shattered by \mathcal{I} , i.e.

Definition 3.2 (VC dimension) *The VC dimension of a range space $\mathcal{R} = (U, \mathcal{I})$, denoted $VCDim(\mathcal{R})$, is*

$$VCDim(\mathcal{R}) = \max\{k : \exists S \subseteq U \text{ such that } |S| = k \text{ and } |\mathcal{I}_S| = 2^k\}.$$

The following combinatorial object, called ε -net, is useful when one wants to find a sample $S \subseteq U$ that intersects every range in \mathcal{I} of a sufficient size.

Definition 3.3 (ε -net) *Let $\mathcal{R} = (U, \mathcal{I})$ be a range space and π be a probability distribution on U . Given $0 < \varepsilon < 1$, a set S is called ε -net w.r.t. \mathcal{R} if*

$$\forall I \in \mathcal{I}, \Pr_{\pi}(I) \geq \varepsilon \implies |I \cap S| \geq 1.$$

When computing ε -nets for a given range space $\mathcal{R} = (U, \mathcal{I})$, we typically build a sample S from elements of U . One can obtain lower bounds for the size of S via standard union bound. However, these bounds usually overestimate $|S|$ since they only take into account the number of points in U or the number of ranges in \mathcal{R} . This issue can be overcome if the VC dimension of the range space that models the problem at hand, denoted k , is finite. The next theorem, proven by Har-Peled and Sharir (2011) [10], states a lower bound for $|S|$ based on k .

Theorem 2 (see [10], Theorem 2.12) *Given $0 < \varepsilon, \delta < 1$, let $\mathcal{R} = (U, \mathcal{I})$ be a range space with $VCDim(\mathcal{R}) \leq k$, let π be a probability distribution on the domain U , and let c be a universal positive constant.*

A collection of elements $S \subseteq U$ sampled w.r.t. π with $|S| = \frac{c}{\varepsilon} (k \ln \frac{1}{\varepsilon} + \ln \frac{1}{\delta})$ is an ε -net with probability at least $1 - \delta$.

As pointed by Löffler and Phillips (2009) [15], c is around $\frac{1}{2}$, but in this paper we leave c as an unspecified constant.

Some of the techniques used in our sampling strategy described in Sections 3.1 and 4 were developed by Riondato and Kornaropoulos (2016) and Riondato and Upfal (2018) [19, 20], where the authors used VC dimension theory, the ε -sample theorem, and Rademacher averages for the estimation of betweenness centrality in a graph. The work of Lima et al. [13, 12] showed how to use sample complexity tools for the estimation of the percolation centrality, which is a generalization of the betweenness centrality. More recently, Cousins et al. (2021) [7] showed improved bounds for the betweenness centrality approximation using Monte–Carlo empirical Rademacher averages, and Lima et al. (2022) [14] used sample complexity tools in the design of a sampling algorithm for the local clustering coefficient of every vertex of a graph.

3.1 Range Space and VC Dimension Results

In this section, we first define the problem in terms of a range space; that is, the problem of computing, with probability at least $1 - \delta$, the shortest paths between the pairs of vertices $(u, v) \in V^2$ that have $c(u, v) \geq \varepsilon$ (Definition 2.2). We show that the VC dimension of the range space that models such problem is constant, which directly impacts in the size of the sample to be used by our algorithms. In fact, we show that this sample size only depends on the parameters of quality and confidence, ε and δ , respectively.

Let $n = |V|$ and \mathcal{T} be the set of n Dijkstra trees of G . Recall that such trees are, by definition, composed by canonical paths. The universe U is defined for the set of all branches of Dijkstra trees, i.e.

$$U = \bigcup_{(a,b) \in V^2: a \neq b} \mathcal{B}_{ab}.$$

For each pair $(u, v) \in V^2$, let p_{uv} be the canonical path from u to v , according to Definition 2.1. Each range τ_{uv} is defined as $\tau_{uv} = \{\mathcal{B}_{ab} \in U : p_{uv} \in S(\mathcal{B}_{ab})\}$. In other words, we can say that \mathcal{B}_{ab} is in the range of (u, v) if \mathcal{B}_{ab} “passes” through a canonical path between u and v . Let $\mathcal{I} = \{\tau_{uv} : (u, v) \in V^2\}$ be the range set. So, $\mathcal{R} = (U, \mathcal{I})$ is the range space defined for our problem.

Now we show how to plug our range space \mathcal{R} with Definition 3.3 so we can use Theorem 2 to bound the sample size that is tight enough for the task that we are tackling. We first show in Theorem 3 that $\Pr_{\pi}(\tau_{uv}) = c(u, v)$. For this result, we have that each tree $T_a \in \mathcal{T}$ is sampled with probability $\pi(T_a) = \frac{1}{n}$ and each branch $\mathcal{B}_{ab} \in T_a$ is sampled with probability $\frac{1}{n-1}$, leading to the probability distribution $\pi(\mathcal{B}_{ab}) = \frac{1}{n(n-1)}$ (which is a proper distribution as the sum is equal to 1). Let $\mathbb{1}_{uv}(\mathcal{B}_{ab})$ be the indicator function that returns 1 if there is some canonical path from u to v as subpath of \mathcal{B}_{ab} , i.e. $\mathcal{B}_{ab} \in \tau_{uv}$, and 0 otherwise. The value of τ_{uv} is equal to counting the individual probabilities of each branch that is in τ_{uv} .

Theorem 3 *For $(u, v) \in V^2$, $\Pr_{\pi}(\tau_{uv}) = c(u, v)$.*

Proof: For fixed $(u, v) \in V^2$ and considering that a branch $\mathcal{B}_{ab} \in U$ is sampled with probability $\pi(\mathcal{B}_{ab}) = \frac{1}{n(n-1)}$, we have

$$\begin{aligned} \Pr_{\pi}(\tau_{uv}) &= \sum_{T_a \in \mathcal{T}} \sum_{\mathcal{B}_{ab} \in T_a} \pi(\mathcal{B}_{ab}) \mathbb{1}_{uv}(\mathcal{B}_{ab}) \\ &= \frac{1}{n(n-1)} \sum_{T_a \in \mathcal{T}} \sum_{\mathcal{B}_{ab} \in T_a} \mathbb{1}_{uv}(\mathcal{B}_{ab}) \\ &= \frac{1}{n(n-1)} \sum_{a \in V} \sum_{b \in V: a \neq b} \mathbb{1}_{uv}(\mathcal{B}_{ab}) \\ &= \frac{t_{uv}}{n(n-1)} = c(u, v). \end{aligned}$$

The first equality follows from the fact that the probability that a branch lies on the range τ_{uv} corresponds to counting the individual probabilities of each branch that is in τ_{uv} . \square

For problems involving shortest paths, such as the ones in [19, 12], it is possible to find a bound for the sample size using VC dimension theory. The referred work typically apply the same proof structure, having a bound based on the vertex-diameter of a graph G , denoted $\text{Diam}_V(G)$, as in Theorem 4 (we present such proof for the sake of completeness). Even though $\text{Diam}_V(G)$ might be as large as n , in particular, this bound is exponentially smaller for graphs with logarithmic vertex-diameter, which may be common in practice. In [19] it is presented a constant VC dimension associated to the set of shortest paths for a graph that contains exactly one shortest path between every pair of vertices. We note that our result generalizes their work, since in our case we do not require such restriction on the input graph.

Although the bound presented in Theorem 4 depends on a combinatorial structure of G , in this work we present an improvement to this result in Theorems 5 and 6, giving a bound that depends only on the desired quality and confidence parameters of the solution. More specifically, for these two theorems we have that $\text{VCDim}(G) = 2$ for a given graph G with respect to a fixed vertex ordering σ , where $\text{VCDim}(G)$ denotes the VC dimension of the range space $\mathcal{R} = (U, \mathcal{I})$ related to a graph G .

Theorem 4 *For a given graph $G = (V, E)$,*

$$\text{VCDim}(G) \leq \lfloor 2 \log \text{Diam}_V(G) + 1 \rfloor.$$

Proof: Let $\text{VCDim}(G) = k$, where $k \in \mathbb{N}$. Then, there is $S \subseteq U$ such that $|S| = k$ and S is shattered by \mathcal{I} . Each $\mathcal{B}_{ab} \in S$ must appear in 2^{k-1} different ranges in \mathcal{I} , from the definition of shattering. On the other hand, \mathcal{B}_{ab} has length at most $\text{Diam}_V(G)$. Then the maximum number of subpaths of \mathcal{B}_{ab} , denoted $|S(\mathcal{B}_{ab})|$, is $\text{Diam}_V(G) \cdot (\text{Diam}_V(G) - 1)$. Thus, the branch \mathcal{B}_{ab} lies in at most $|S(\mathcal{B}_{ab})|$ ranges, and therefore,

$$2^{k-1} \leq |S(\mathcal{B}_{ab})| \leq \text{Diam}_V(G) \cdot (\text{Diam}_V(G) - 1) \leq \text{Diam}_V(G)^2.$$

Solving for k , $\text{VCDim}(G) = k \leq \lfloor 2 \log \text{Diam}_V(G) + 1 \rfloor$. \square

For Theorems 5 and 6, we introduce the definition of *meeting path* between two canonical paths P_1 and P_2 , and in Lemma 5 we prove that there is only one such path between P_1 and P_2 . We use this fact to prove that $\text{VCDim}(G) \leq 2$ in Theorem 5. On this theorem, we show that a set of

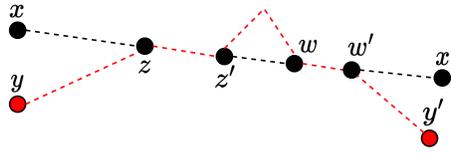


Figure 3: Illustration of the situation described in Lemma 5, where the paths P_1 and P_2 are depicted in black and red colors, respectively.

three paths $S = \{P_1, P_2, P_3\}$ cannot be shattered by the range set. More specifically and without losing of generality, we show that any meeting path between P_1 and P_3 is also a meeting path of P_2 , and therefore, all the possible subsets of S but $\{P_1, P_3\}$ can be found in the intersection of S with the range set.

Definition 3.4 Consider two different canonical paths P_1 and P_2 . We say that a canonical path $Z = (z, \dots, z')$ is a meeting path between P_1 and P_2 if Z is a maximal (z, z') -subpath of P_1 and P_2 .

On Definition 3.4, we consider two paths intersecting if they share at least one edge.

Lemma 5 Consider two different canonical paths P_1 and P_2 . Let Z be a meeting path between P_1 and P_2 . Then Z is the only meeting path between both paths in G .

Proof: Let $P_1 = (x, \dots, x')$, $P_2 = (y, \dots, y')$, and $Z = (z, \dots, z')$. Suppose that Z is a meeting path between P_1 and P_2 and suppose that it is not unique. Let $W = (w, \dots, w')$ be another meeting path between P_1 and P_2 in G (Figure 3). Note that Z and W are disjoint, otherwise the concatenation of both paths would contradict the maximality of Z and W . Without loss of generality, we may assume the following:

- Z is contained in the (x, z') -subpath of P_1 and in the (y, z') -subpath of P_2 , with z' closer to x and to y than to x' and y' in P_1 and P_2 , respectively;
- W is contained in the (w, x') -subpath of P_1 and in the (w, y') -subpath of P_2 , with w closer to x' and to y' than to x and y in P_1 and P_2 , respectively.

Let D be the CP from z' to w in G . Since P_1 and P_2 are canonical paths, by Theorem 1, the (z', w) -subpath of P_1 and the (z', w) -subpath of P_2 must be equal do D . Let Z' be the concatenation of Z , D , and W . Then Z' is a meeting path between P_1 and P_2 that contradicts the maximality of Z . □

Theorem 5 For a given graph $G = (V, E)$ and a fixed ordering σ over V ,

$$VCDim(G) \leq 2.$$

Proof: Suppose that $VCDim(G) > 2$. Then, from the definition of VC dimension, there is a set of canonical paths $\mathcal{S} = \{P_1, P_2, P_3\}$ that is shattered by \mathcal{I} . These paths are described as $P_1 = \{u, \dots, v\}$, $P_2 = \{u', \dots, v'\}$, and $P_3 = \{u'', \dots, v''\}$. Let W be the (w, w') -subpath of P_1 that is also contained in P_2 and P_3 . From the definition of shattering, this path must exist so that $\tau_{ww'} \cap \mathcal{S} = \{P_1, P_2, P_3\}$. Let x be the farthest predecessor of w in P_1 such that, w.l.o.g.,

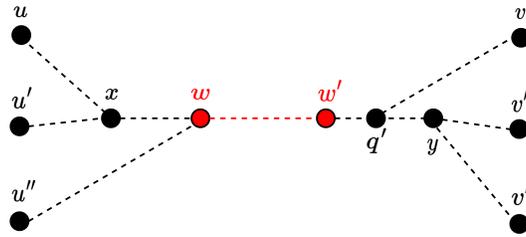


Figure 4: Illustration of paths P_1 and P_2 having a (x, w) -path as a meeting path, and of paths P_2 and P_3 having a (q', y) -path as a meeting path.

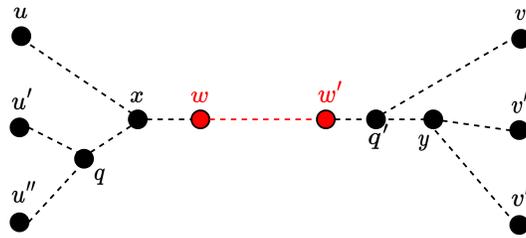


Figure 5: Case where $\tau_{xw} \cap S = \{P_1, P_2, P_3\}$, for $P_1 = (u, \dots, v)$, $P_2 = (u', \dots, v')$, and $P_3 = (u'', \dots, v'')$.

the (x, w) -subpath of P_1 , denoted X , is also contained in P_2 but not in P_3 . Let y be the farthest successor of w in P_1 such that a (q', y) -subpath of P_1 , denoted Y , is also contained in P_3 but not in P_2 . Note that X and Y must exist so that $\tau_{xw} \cap S = \{P_1, P_2\}$ and $\tau_{q'y} \cap S = \{P_1, P_3\}$ (Figure 4).

We now show that either P_3 can only have a meeting path with P_2 and not with P_1 (and vice-versa), or it can only have a meeting path that is common to P_1 and P_2 at the same time (Figures 5, 6, and 7). Suppose that there is a (q, x) -subpath of P_2 that is contained in P_3 but not in P_1 , as depicted in Figure 5. The vertex q is not contained in X , but P_2 and P_3 must pass through W . Besides, the CP from u' to v' is not the same as the one from v' to u' (and correspondingly for u'' and v''), so a meeting path between P_3 and a shortest path that goes from v' to u' is not the same meeting path between P_2 and P_3 . Therefore, the (q, x) -subpath pass through the (x, w) -subpath of P_2 . From Lemma 5, all the vertices from q to w' must be the same in P_2 and P_3 . Hence, P_3 goes through x , and from our initial assumption, P_2 does not have any intersection with a vertex that comes before x in P_1 . Besides, P_3 goes through q' and Y . Therefore, any subpath of P_2 starting in q is also a subpath of P_3 . This contradicts that $\tau_{xw} \cap S = \{P_1, P_2\}$ since $\tau_{xw} \cap S = \{P_1, P_2, P_3\}$.

Consider now the (q', v) -subpath of P_1 , denoted P'_1 . Suppose that P_3 has an intersection with P_1 on a (r, r') -subpath of P'_1 . Such path is depicted in red in Figure 6. Note that in fact this cannot happen, otherwise G would have more than one canonical shortest path from q' to r in G .

Now, consider the (q', v') -subpath of P_2 , denoted P'_2 . Suppose that P_3 has an intersection with a (r, r') -subpath of P'_2 (Figure 7). From our initial assumption, P_3 goes through W and Y , so it passes through q' , and q' reaches r . Hence, from Lemma 5, all the vertices from q' to r' must be the same in P_2 and P_3 . In this case, P_3 does not contain a (r', w) -subpath, otherwise P_1 and

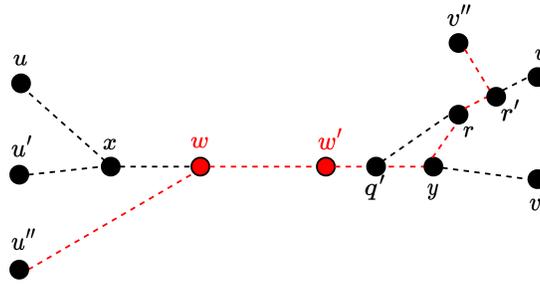


Figure 6: Case of an intersection of the paths P_1 and P_3 that is prohibited. In this case, there is more than one canonical shortest path from q' to r in G .

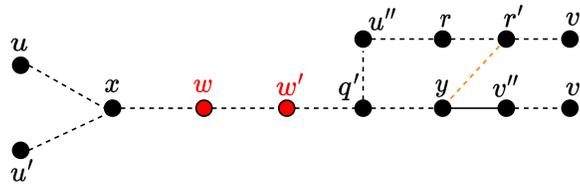


Figure 7: Case where $\tau_{q'y} \cap S = \{P_1\}$, for $P_1 = (u, \dots, v)$, $P_2 = (u', \dots, v')$, and $P_3 = (u'', \dots, v'')$. The orange dashed path correspond to a shortest path that cannot happen.

P_3 would form a cycle starting and ending in r' . Besides, P_3 does not have a (r', y) -subpath or a (r', y') -subpath, for any $y' \in \text{Inn}(Y)$, otherwise that would be two different CPs from r' to y' . Hence, P_3 does not pass through the (q', y) -subpath of P_1 , contradicting that $\tau_{q'y} \cap S = \{P_1, P_3\}$ since $\tau_{q'y} \cap S = \{P_1\}$. □

Theorem 6 For a given graph $G = (V, E)$ and a fixed ordering σ over V ,

$$VCDim(G) \geq 2.$$

Proof: Consider the graph as in Figure 8. Then, for $P_1 = (a, b, c, d, e, f)$, $P_2 = (g, b, c, d, e, h)$, and $S = \{P_1, P_2\}$, we have: $\tau_{ac} = \{P_1\}$, $\tau_{gc} = \{P_2\}$, $\tau_{cd} = \{P_1, P_2\}$, and $\tau_{aj} = \emptyset$.

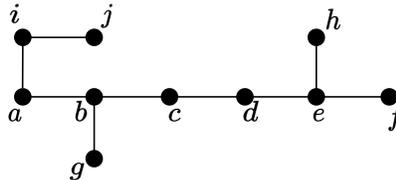


Figure 8: Graph with $VCDim(G) \geq 2$. □

4 Algorithms

For an undirected graph $G = (V, E)$ with non-negative edges weights, with $n = |V|$ and $m = |E|$, we first present in Section 4.1 a modified version of Dijkstra's algorithm which takes into consideration a given vertex ordering σ , and then we show that the shortest paths in the SPT computed by the algorithm are canonical paths. Then, in Section 4.2 we present an algorithm for the relaxed APSP problem that returns, with probability at least $1 - \delta$, the shortest paths that appear as subpaths of at least a proportion ε of all canonical paths.

4.1 Modified Dijkstra

In this section we present a modification of Dijkstra's algorithm presented in [6]. Dijkstra's algorithm, for a given vertex s , outputs a SPT, denoted T_s , rooted in s . This algorithm maintains in every step a set S such that every vertex in S has its distance from s already computed. At every step, a vertex v in $V \setminus S$ with minimum estimated distance from s is selected to be added in S . An edge $(u, w) \in E$ is *relaxed* if the minimum distance from s to u plus the weight of (u, w) improves the minimum distance from s to w .

The main difference between the modified algorithm that we present here and the original one is the tie-breaking criterion for the selection of edges to be added in a shortest path. In a given step of the modified Dijkstra, if there are multiple vertices in $V \setminus S$ with the same estimation for minimum distance from s , then the one with minimum index in σ is chosen to be added in S . Additionally, let u be a vertex that has been just inserted in T_s in a given iteration. For every neighbor y of u in $V \setminus S$ for which the algorithm relaxes the edge (u, y) , the ordering is taken into consideration so that if $d(s, u) + \omega(u, y) = d(s, u') + \omega(u', y)$, for some u' in S , then the tie-breaking for the shortest (s, y) -path depends on which vertex between u and u' has the minimum index in σ .

Theorem 7 shows that the modified Dijkstra's algorithm correctly computes all the canonical paths from a source s to any other vertex in V with respect to σ . Note that S is a priority queue that is also modified to give higher priority to vertices with lowest indexes in σ in the case of ties in the vertices selection. We observe, however, that these modifications do not increase the running time of the priority queue operations. In particular, when a vertex is chosen to be included in a shortest path and the priority queue needs to be rearranged, the maximum number of comparisons between the vertices performed in this task is at most the number of comparisons made by the original Dijkstra algorithm.

Theorem 7 *All shortest paths computed by a modified Dijkstra's algorithm with respect to a given vertex ordering σ are canonical paths.*

Proof: (Sketch) Similar to the proof of correctness of the original Dijkstra's algorithm presented in [6] (Theorem 22.6), the proof is by induction on the size of S .

Let s be the source vertex. For each $u \in V$, let $\tilde{d}(s, u)$ the estimated minimum distance from s to u in a given step of the algorithm. For $|S| = 0$, the set S is empty and then this base is trivially true. For the base where $|S| = 1$, we have $S = \{s\}$, and then $\tilde{d}(s, s) = d(s, s) = 0$. Besides, s does not have a predecessor, since it is the source, so the base is also true for this case. For the inductive step, we have the following hypothesis: for all $v \in S$, we have that $\tilde{d}(s, v) = d(s, v)$ and the predecessor of v in the Dijkstra tree of s is the one with minimum index in σ . Proving that $\tilde{d}(s, v) = d(s, v)$ follow the same arguments of the proof of correctness in [6] for the original Dijkstra's algorithm.

In order to prove that the predecessor of v in the Dijkstra tree of s , denoted v' , is the one with minimum index in σ among all possible predecessors of v , we prove that all edges (z, v) where $\tilde{d}(s, v) = \tilde{d}(s, z) + \omega(z, v)$ were examined when the edge (v', v) were relaxed. Consider, by contradiction, that there is some vertex u' that has the minimum index in σ among all possible predecessors of v , but that the edge (u', v) was not examined before vertex v is added to S . If the edge (u', v) was not examined, then v was added in S before u' . In this case, this happened either because $\tilde{d}(s, v) < \tilde{d}(s, u')$ or because $\tilde{d}(s, v) = \tilde{d}(s, u')$ and $\sigma(v) < \sigma(u')$. However, in both cases, then u' could not be the predecessor of v , since $\tilde{d}(s, u')$ should be strictly smaller than $\tilde{d}(s, v)$ to be considered as a possible predecessor of v . Hence, all $y \in S$ with $\tilde{d}(s, y) < \tilde{d}(s, v)$ should have been examined before v , and hence, v' is the predecessor of v with minimum index in σ among all such vertices. This value never changes again once v is added in S . \square

4.2 Computing Shortest Paths

Given $0 < \varepsilon, \delta < 1$, Algorithm 1 computes, with probability $1 - \delta$, the distances between pairs of vertices $(u, v) \in V^2$ such that $c(u, v) \geq \varepsilon$. We also briefly describe the necessary modifications on the algorithm so that the shortest path associated to such distances be also computed.

Algorithm 1: PROBABILISTICALLPAIRSHORTESTPATHS(G, ε, δ)

```

input : non-negative weighted graph  $G = (V, E)$  with  $n = |V|$ , parameters  $0 < \varepsilon, \delta < 1$ .
output: distance  $d_{uv}$ , for each  $(u, v) \in V^2$  s.t.  $c(u, v) > \varepsilon$ , with probability  $1 - \delta$ .
1 for  $i \leftarrow 1$  to  $\lceil \frac{\varepsilon}{\delta} (2 \ln \frac{1}{\varepsilon} + \ln \frac{1}{\delta}) \rceil$  do
2   sample  $a \in V$  with probability  $1/n$ 
3    $T_a \leftarrow$  SINGLESOURCESHORTESTPATHS( $a$ )           /* modified Dijkstra */
4   sample  $b \in V \setminus \{a\}$  with probability  $1/(n - 1)$ 
5    $\mathcal{B}_{ab} \leftarrow$  shortest path from  $a$  to  $b$  in  $T_a$ 
6   for each  $(u, v) \in \mathcal{B}_{ab} \times \mathcal{B}_{ab}$  do           /*  $u$  closer to  $a$ ,  $v$  closer to  $b$  */
7      $d_{uv} \leftarrow d_{av} - d_{au}$                  /*  $d_{au}$  and  $d_{av}$  come from  $T_a$  */
8 return each  $d_{uv}$  in the distances table

```

Theorem 8 Consider a (u, v) -path such that $c(u, v) \geq \varepsilon$. Algorithm 1 computes the exact distance between u and v with probability $1 - \delta$.

Proof: Algorithm 1 samples several branches and we first assume that such samples are an ε -net (we show later that this is indeed true). Recalling the range space modeling (Section 4.2), the sample of branches is denoted by S and the (u, v) -path is related to a range τ_{uv} .

As, by lines 2 and 4, the branch is sampled with probability $1/n(n - 1)$ then, by Theorem 3, we have that $c(u, v) = \Pr(\tau_{uv})$. Thus, as $c(u, v) \geq \varepsilon$, so $\Pr(\tau_{uv}) \geq \varepsilon$. As we are assuming that the sample is an ε -net, by Definition 3.3, then $|\tau_{uv} \cap S| \geq 1$ for all τ_{uv} such that $\Pr(\tau_{uv}) \geq \varepsilon$. That is, since $c(u, v) \geq \varepsilon$ then at least one branch of the sample S contains the (u, v) -path. If a branch \mathcal{B}_{ab} in S contains the (u, v) -path, then in line 3 the exact distance between u and v is computed, since the (u, v) -path which is a subpath of the shortest path from a to b is also minimal, so its distance d_{uv} can be computed as $d_{av} - d_{au}$.

Now it remains to prove that the sample S is indeed an ε -net. Note that in lines 1–7, the loop is executed $k = \lceil \frac{\varepsilon}{\delta} (2 \ln \frac{1}{\varepsilon} + \ln \frac{1}{\delta}) \rceil$ times, so our sample has at least size k . By Theorems 2, 5, and 6, this sample size is sufficient for it to be an ε -net with probability at least $1 - \delta$. \square

Theorem 9 *Algorithm 1 has running time $\mathcal{O}(m + n \log n + (\text{Diam}_V(G))^2)$.*

Proof: Lines 2 and 4 take constant time and line 5 takes linear time. Line 3 (the modified Dijkstra) runs in $\mathcal{O}(m + n \log n)$, as the modifications do not change the running time of the original Dijkstra’s algorithm. The loop in line 6 takes time $\mathcal{O}((\text{Diam}_V(G))^2)$ since the length of \mathcal{B}_{ab} cannot be greater than the vertex diameter of the graph. The distances returned by Dijkstra’s algorithm in line 3 are stored in a table d . Since operations of insertion, deletion, and search on this data structure take time $\mathcal{O}(1)$, then updating table d takes time $\mathcal{O}(1)$. Assuming that ε and δ are constants, the number of loop iterations in lines 1–7 is constant, and the result follows. \square

As it is common to APSP and search algorithms, Algorithm 1 also constructs a data structure from which, for all vertices (u, w) , a shortest path from u to w can be retrieved. We can store the predecessors of each vertex that is in \mathcal{B}_{ab} so that a (u, v) -subpath of \mathcal{B}_{ab} can be retrieved by a backward traversing from v to u on these predecessors. This modification does not change the execution time of the original algorithm.

In the remainder of this section we are interested in determining the smallest value of ε for which our algorithm would still perform on strictly subcubic time. For this, we drop the assumption that ε is constant and therefore write it as a function of n , denoted by $\varepsilon(n)$.

Let k be the sample size (which impacts on the number of times line 1 of Algorithm 1 is executed). Then $k = \mathcal{O}\left(\frac{1}{\varepsilon(n)} \ln \frac{1}{\varepsilon(n)}\right)$, and the running time of Algorithm 1 becomes $\mathcal{O}(k \cdot (m + n \log n + (\text{Diam}_V(G))^2))$. In the worst case $m = \mathcal{O}(n^2)$ and then its running time is $\mathcal{O}(k \cdot n^2)$. As the best conjectured time is $\mathcal{O}(n^{3-c})$, for a constant $c > 0$ [28], then we are looking for the value of $\varepsilon(n)$ such that the time of our algorithm is upper bounded by $\mathcal{O}(n^{3-c})$, i.e. $\mathcal{O}(k \cdot n^2) = \mathcal{O}(n^{3-c})$. Thus $k = n^{1-c}$, i.e.

$$\frac{1}{\varepsilon(n)} \ln \frac{1}{\varepsilon(n)} = n^{1-c}.$$

Solving for $\varepsilon(n)$, we have $\varepsilon(n) = \frac{W_0(n^{1-c})}{n^{1-c}}$, where $W_0(n^{1-c})$ is the branch 0 of the Lambert-W function [26]. To simplify the notation, let $n' = n^{1-c}$. If $n' \geq e$, then a known bound [11] for $W_0(n')$ is $W_0(n') = \ln n' - \ln \ln n' + \Theta\left(\frac{\ln \ln n'}{\ln n'}\right)$. Therefore $\varepsilon(n) = \frac{\ln n' - \ln \ln n' + \Theta\left(\frac{\ln \ln n'}{\ln n'}\right)}{n'}$.

Note that the smallest value for $c(u, v)$, for a pair $(u, v) \in V^2$, is $1/n(n - 1)$, which is the case for a path that is not strictly contained in any other path. So, to compute the distance of paths with such small value, we have to use ε so small that the execution time exceeds that of the best existing algorithms [28, 18]. Nevertheless, by the reasoning above, we note that we can set ε as small as $\Theta\left(\frac{\ln n'}{n'}\right)$.

5 Concluding Remarks

In this paper we present a range space having the domain composed by the shortest paths of a graph G where there is one shortest path for each pair of vertices in G . We show that the VC dimension of such range space is 2. We show that this result can be applied to bound the sample size required for an approximation algorithm for a relaxed version of the All-pairs shortest path problem (APSP). In this version, we consider the set S , which is a subset of all shortest paths from a graph G such that S contains exactly one shortest path between every pair of distinct vertices of G . We compute, with probability at least $1 - \delta$, the shortest paths of G that appear as subpath of at least a proportion ε of all shortest paths in the set S , for $0 < \varepsilon, \delta < 1$. We

present a $\mathcal{O}(m + n \log n + (\text{Diam}_V(G))^2)$ running time algorithm for this task. We show that a sample of shortest paths of size $\lceil \frac{\varepsilon}{\delta} (2 \ln \frac{1}{\varepsilon} + \ln \frac{1}{\delta}) \rceil$ is sufficient for achieving the desired result. So, in an application where one might be interested only in computing “important” shortest paths the algorithm is rather efficient and it depends only on the parameters ε and δ (classical approaches in literature based in union bound, for example, typical require sample sizes that depend on the size of the input).

An open question that we are particularly interested is the connection between ε and n or $\text{Diam}_V(G)$ for specific input distributions. For the general case, trivially setting $\varepsilon = \frac{1}{n(n-1)}$, we have a guarantee that every shortest path in G is computed with probability $1 - \delta$, but that would yield an algorithm with running time exceeding $\mathcal{O}(n^3)$. This may not be a surprise since APSP may not admit a strictly subcubic algorithm. Nevertheless, we show that if ε is at least $\frac{\ln n' - \ln \ln n' + \Theta\left(\frac{\ln \ln n'}{\ln n'}\right)}{n'}$, where $n' = n^{1-c}$, the running time of our algorithm is $\mathcal{O}(n^{3-c})$, for $c > 0$.

Instead of fixing a vertex ordering and using the strategy of canonical paths, one may ask whether a simpler strategy would not be enough, such as adding some perturbation on the edges weights for artificially changing the input graph so that the shortest paths are unique. However, rather than changing the input graph, we prefer to assume that we are completely subordinated to the input distribution, which seems to be reasonable for a sampling algorithm. In fact, one can find pathological inputs where different vertex orderings can generate significantly different values of $c(u, v)$, for $(u, v) \in V^2$, but since we are dealing with sampling, that is not a serious issue.

An extensive experimental evaluation is out of the scope of our work, which is theoretical in nature, however, we observe that, in practice, the ordering of vertices may have little impact in the value of $c(u, v)$. Some preliminary experiments¹ showed that, as expected, both the average standard deviation and the maximum standard deviation are small for the values of $c(u, v)$ (average standard deviation is 10^{-6} and maximum standard deviation is 10^{-3}). However, one might be interested in investigating more rigorously this question, e.g., specifying an input distribution and performing a probabilistic analysis.

Acknowledgements

This work was partially funded by the Coordination for the Improvement of Higher Education Personnel (CAPES) and the National Council for Scientific and Technological Development (CNPq).

References

- [1] A. Abboud and V. Williams. Popular conjectures imply strong lower bounds for dynamic problems. In *2014 IEEE 55th Annual Symposium on Foundations of Computer Science*, pages 434–443, 2014. doi:10.1109/FOCS.2014.53.
- [2] A. Abboud, V. Williams, and H. Yu. Matching triangles and basing hardness on an extremely popular conjecture. *SIAM Journal on Computing*, 47(3):1098–1122, 2018. doi:10.1145/2746539.2746594.

¹We performed a preliminary battery of experiments using real-world graphs from Network Repository (available online in <https://networkrepository.com/networks.php>). We used 10 graphs with vertex set size ranging from 200 to 9500 (we included sparse and dense graphs).

- [3] M. Anthony and P. L. Bartlett. *Neural Network Learning: Theoretical Foundations*. Cambridge University Press, New York, 1st edition, 2009.
- [4] A. Brodnik and M. Grgurovič. Solving all-pairs shortest path by single-source computations: Theory and practice. *Discrete applied mathematics*, 231:119–130, 2017. doi:10.1016/j.dam.2017.03.008.
- [5] T. M. Chan. All-pairs shortest paths for unweighted undirected graphs in $o(mn)$ time. *ACM Trans. Algorithms*, 8(4), 2012. doi:10.1145/2344422.2344424.
- [6] T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein. *Introduction to algorithms*. MIT press, 2022.
- [7] C. Cousins, C. Wohlgemuth, and M. Riondato. Bavarian: Betweenness centrality approximation with variance-aware rademacher averages. In *27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, pages 196–206, 2021. doi:10.1145/3577021.
- [8] D. Dor, S. Halperin, and U. Zwick. All-pairs almost shortest paths. *SIAM Journal on Computing*, 29(5):1740–1759, 2000. doi:10.1137/S0097539797327908.
- [9] P. Eirinakis, M. D. Williamson, and K. Subramani. On the shoshan-zwick algorithm for the all-pairs shortest path problem. *Journal of Graph Algorithms and Applications*, 21(2):177–181, 2017. doi:10.7155/jgaa.00410.
- [10] S. Har-Peled and M. Sharir. Relative (p, ϵ) -approximations in geometry. *Discrete & Computational Geometry*, 45(3):462–496, 2011. doi:10.1007/s00454-010-9248-1.
- [11] A. Hoorfar and M. Hassani. Inequalities on the lambert w function and hyperpower function. *J. Inequal. Pure and Appl. Math*, 9(2):5–9, 2008.
- [12] A. M. Lima, M. V. da Silva, and A. Vignatti. Percolation centrality via rademacher complexity. *Discrete Applied Mathematics*, 2021. doi:10.1016/j.dam.2021.07.023.
- [13] A. M. Lima, M. V. da Silva, and A. L. Vignatti. Estimating the percolation centrality of large networks through pseudo-dimension theory. In *26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 1839–1847, 2020. doi:10.1145/3394486.3403235.
- [14] A. M. Lima, M. V. G. da Silva, and A. L. Vignatti. Estimating the clustering coefficient using sample complexity analysis. In *LATIN 2022: Theoretical Informatics*, pages 328–341, Cham, 2022. Springer International Publishing. doi:10.1145/2930889.2930939.
- [15] M. Löffler and J. M. Phillips. Shape fitting on point sets with probability distributions. In *European symposium on algorithms*, pages 313–324. Springer, 2009. doi:10.1007/978-3-642-04128-0_29.
- [16] M. Mitzenmacher and E. Upfal. *Probability and Computing: Randomization and Probabilistic Techniques in Algorithms and Data Analysis*. Cambridge University Press, New York, 2nd edition, 2017.
- [17] M. Mohri, A. Rostamizadeh, and A. Talwalkar. *Foundations of Machine Learning*. The MIT Press, Cambridge, 2012.

- [18] S. Pettie and V. Ramachandran. Computing shortest paths with comparisons and additions. In *13th Annual ACM-SIAM Symposium on Discrete Algorithms*, SODA '02, pages 267–276, Philadelphia, PA, USA, 2002. Society for Industrial and Applied Mathematics. doi:10.5555/545381.545417.
- [19] M. Riondato and E. M. Kornaropoulos. Fast approximation of betweenness centrality through sampling. *Data Mining and Knowledge Discovery*, 30(2):438–475, 2016. doi:10.1145/2556195.2556224.
- [20] M. Riondato and E. Upfal. Abra: Approximating betweenness centrality in static and dynamic graphs with rademacher averages. *ACM Trans. Knowl. Discov. Data*, 12(5):61:1–61:38, 2018. doi:10.1145/3208351.
- [21] L. Roditty and A. Shapira. All-pairs shortest paths with a sublinear additive error. *ACM Transactions on Algorithms (TALG)*, 7(4):1–12, 2011. doi:10.1145/2000807.2000813.
- [22] L. Roditty and V. Vassilevska Williams. Fast approximation algorithms for the diameter and radius of sparse graphs. In *45th annual ACM symposium on Theory of computing*, pages 515–524, 2013. doi:10.1145/2488608.2488673.
- [23] S. Shalev-Shwartz and S. Ben-David. *Understanding Machine Learning: From Theory to Algorithms*. Cambridge University Press, New York, 2014.
- [24] A. Shoshan and U. Zwick. All pairs shortest paths in undirected graphs with integer weights. In *40th Annual Symposium on Foundations of Computer Science (Cat. No. 99CB37039)*, pages 605–614, 1999. doi:10.1109/SFFCS.1999.814635.
- [25] V. Vassilevska Williams. Hardness of easy problems: Basing hardness on popular conjectures such as the strong exponential time hypothesis (invited talk). In *10th International Symposium on Parameterized and Exact Computation (IPEC 2015)*. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2015. doi:10.4230/LIPIcs.IPEC.2015.17.
- [26] E. Weisstein. Lambert w-function, from wolfram math-world, 2013.
- [27] R. Williams. Faster all-pairs shortest paths via circuit complexity. In *46-th Annual ACM Symposium on Theory of Computing*, STOC'14, pages 664–673, New York, 2014. ACM. doi:10.1145/2591796.2591811.
- [28] R. Williams. Faster all-pairs shortest paths via circuit complexity. *SIAM Journal on Computing*, 47(5):1965–1985, 2018. doi:10.1145/2591796.2591811.