

Cut Locus and Topology from Surface Point Data

Tamal K. Dey* Kuiyu Li†

March 24, 2009

Abstract

A cut locus of a point p in a compact Riemannian manifold M is defined as the set of points where *minimizing* geodesics issued from p stop being minimizing. It is known that a cut locus contains most of the topological information of M . Our goal is to utilize this property of cut loci to decipher the topology of M from a point sample. Recently it has been shown that Rips complexes can be built from a point sample P of M systematically to compute the Betti numbers, the rank of the homology groups of M . Rips complexes can be computed easily and therefore are favored over others such as restricted Delaunay, alpha, Čech, and witness complex. However, the sizes of the Rips complexes tend to be large. Since the dimension of a cut locus is lower than that of the manifold M , a subsample of P approximating the cut locus is usually much smaller in size and hence admits a relatively smaller Rips complex.

In this paper we explore the above approach for point data sampled from surfaces embedded in any high dimensional Euclidean space. We present an algorithm that computes a subsample P' of a sample P of a 2-manifold where P' approximates a cut locus. Empirical results show that the first Betti number of M can be computed from the Rips complexes built on these subsamples. The sizes of these Rips complexes are much smaller than the one built on the original sample of M .

*Department of Computer Science and Engineering, The Ohio State University, Columbus, OH 43210, USA. Email: tamaldey@cse.ohio-state.edu

†Department of Computer Science and Engineering, The Ohio State University, Columbus, OH 43210, USA. Email: liku@cse.ohio-state.edu

1 Introduction

A considerable amount of interest has been generated recently in applying geometric and topological techniques to data analysis in high dimensional spaces [1, 7, 10, 11, 20, 24, 26, 27]. Assuming that the data is sampled from a low dimensional manifold lying in a high dimensional space, the results in these works facilitate algorithms that ‘learn’ different properties of the manifold. We are specifically interested in extracting the topology (information about homology) of the manifold from its point data.

Recently a few algorithms have been proposed for the problem which have theoretical guarantees [3, 5, 10, 24]. These algorithms are theoretically sound but are not practical. They dwell on data structures such as Delaunay triangulations and alpha shapes that have impractically high computational cost in large dimensions. Alternative data structures such as witness complex, Čech complex, and Rips complexes have been proposed [7, 25] to counter this problem. Rips complexes can be computed more easily than the others and so become an attractive choice [8, 18] in applications. Taking this view point, Chazal and Oudot [7] show how one can build a hierarchy of Rips complexes from a point cloud data and then use topological persistence [15, 27] to compute the Betti numbers of the sampled manifold. However, the size of a Rips complex is relatively large and that becomes a bottleneck for computing persistent Betti numbers from them. It is this consideration which motivates our work.

We utilize a well known structure called *cut locus* in differential geometry to cut down the size of the Rips complexes. Let p be any point in a m -dimensional smooth compact Riemannian manifold M . The cut locus $C(p) \subset M$ is the space of points where the minimizing geodesics issued from p stop being minimizing. It is known that $M \setminus C(p)$ is a ball and hence most of the topology of M is contained in $C(p)$. Specifically, ranks of all homology groups (under \mathbb{Z}_2 coefficient ring) of $C(p)$ coincide with those of M except for the full dimensional one. The cut locus being one dimensional lower object than M can be approximated by a subsample of size much smaller than a sample of M . As a result the Rips complexes get much smaller which in turn facilitate computations of persistent homology groups, see Table 1.

In this paper we explore the above approach for surfaces embedded in an arbitrary Euclidean space, that is, M is a compact smooth 2-manifold sitting in \mathbb{R}^k for some $k > 2$. We assume M to be connected and hence only its one dimensional homology group is interesting. We present an algorithm that computes a subsample $P' \subset P$ from a sample P of M where P' approximates a cut locus $C(p)$ when P is sufficiently dense. We distinguish our set up from the framework where M is presented with some linear approximation. Cut loci in the presence of an explicit representation of the surface have been used to compute various types of optimal cycles on the surface, see [12, 16]. One may argue that a linear approximation of the surface can be computed from its point sample first, and then known methods for computing a cut locus can be used. Since we are considering M sitting in high dimensional embedding space, this option is not very practical though theoretically possible. Also, our ulterior goal is to explore the cut locus approach for general dimensional manifold. This paper is a step toward that goal.

2 Geodesics and cut locus

We briefly review some of the key concepts related to geodesics, see [14] for details. Let $M \subset \mathbb{R}^k$ be a compact, connected, smooth manifold without boundary. Assume that the metric in M is induced by the scalar product $\langle \cdot, \cdot \rangle$ in \mathbb{R}^k .

Geodesics. A curve $\gamma: I \subset \mathbb{R} \rightarrow M$ is a *geodesic* if the acceleration representing the rate of change of the tangent $\dot{\gamma}(t)$ has no component along M for all $t \in I$. More formally, the covariant derivative $\frac{D}{dt}(\dot{\gamma}(t))$ is 0 for all $t \in I$. Given a vector u in the tangent space TM_p at a point $p \in M$, there is a geodesic $\gamma(t)$ parameterized by arc lengths where $\gamma(0) = p$ and $\dot{\gamma}(0) = u/\|u\|$. The geodesic γ is said to be *issued* from p . Notice that two points p and q in M may have multiple geodesics between them. Among them, the ones minimizing the length (if they exist) are called the *minimizing geodesics* between p and q . Since M is compact, it is geodesically complete, implying that any two points admit a minimizing geodesic. If the minimizing geodesic between $p, q \in M$ is unique, we denote it as γ_{pq} with the understanding that $\gamma_{pq}(0) = p$.

Distances. One can define the distance of a point p to a set $X \subseteq M$ as $d_M(p, X) = \inf_{x \in X} \ell_{px}$ where ℓ_{px} is the length of a minimizing geodesic between p and x in M . We use similar notation $d_E(p, X)$ to denote the Euclidean distance between a point p and a subset X of \mathbb{R}^k . Abusing the notation we write $d_M(p, q) = d_M(p, \{q\})$ and $d_E(p, q) = d_E(p, \{q\})$ for any two points p and q . It is known that $d_E(p, q) \leq d_M(p, q)$ where $p, q \in M \subset \mathbb{R}^k$. We

also have Hausdorff distances d_H^E and d_H^M between two sets X and Y defined as

$$\begin{aligned} d_H^E(X, Y) &= \max\{\sup_{y \in Y} d_E(y, X), \sup_{x \in X} d_E(x, Y)\} \\ d_H^M(X, Y) &= \max\{\sup_{y \in Y} d_M(y, X), \sup_{x \in X} d_M(x, Y)\}. \end{aligned}$$

Exponential map. Let $p \in M$ be an arbitrary point. We are interested in the geodesics issued from p . There is a natural map called the *exponential map* which takes a vector in the tangent space TM_p at p and maps it to a point on the geodesic issued from p by going over a distance of the length $\|u\|$. Formally, $\exp_p: TM_p \rightarrow M$ where $\exp_p(u) = \gamma(\|u\|)$ so that $\gamma(0) = p$ and $\dot{\gamma}(0) = \frac{u}{\|u\|}$. Since M is compact, the map \exp_p is defined for entire TM_p meaning that each geodesic issued from p continues to be geodesic for the infinite interval $[0, \infty]$. However, such a geodesic may cease to be minimizing at some point.

Cut point and locus. A *cut point* of a geodesic γ issued from p is the point where γ ceases to be minimizing. The locus $C(p)$ of all cut points on geodesics issued from p is called the *cut locus* of p , see Figure 1. There is a related concept called *conjugate locus*. This is the locus of all *conjugate points* where the exponential map is critical. Formally, a point $q = \gamma(t)$ is a conjugate point of $p = \gamma(0)$ if and only if $t\dot{\gamma}(0)$ is a critical point of \exp_p .

At a cut point $q \in C(p)$ of a geodesic γ , only two things may happen:

- (a) Either there is another minimizing geodesic σ starting from p so that $\sigma(t) = \gamma(t) = q$ for some $t \in (0, \infty]$, or
- (b) q is the first conjugate point of p along γ .

In Figure 1, point q satisfies (a) and point s_0 satisfies (b). It is obvious that, for any point $q \in M \setminus C(p)$, the geodesic between p and q which has not crossed $C(p)$ is minimizing.

Injectivity radius and reach. For an m -dimensional manifold M , the exponential map allows us to map rays from the tangent space $TM_p \approx \mathbb{R}^m$ to the geodesics in M . Denote an open Euclidean ball with center at $0 = \exp_p^{-1}(p)$ and radius r as $B(0, r)$. The map \exp_p is injective on $B(0, r)$ if and only if r is smaller than or equal to the geodesic distance of p to $C(p)$. This motivates the definition of *injectivity radius* of M given by

$$i(M) = \inf_{p \in M} d_M(p, C(p)).$$

Injectivity radius can be seen as the intrinsic counterpart of a well known extrinsic measure called the *reach*,

$$\rho(M) = \inf_{p \in M} d_E(p, Y)$$

where Y is the medial axis of M [17]. Because of the property (b) of the cut points, \exp_p on $B(0, r)$ is not only injective but also a diffeomorphism if $r < i(M)$. The image $\exp_p(B(0, r))$ is a geodesic ball of radius r in M centered around p .

Cut locus topology. One may observe that the injectivity of \exp_p can be extended to the entire open set $M \setminus C(p)$. It follows that $M \setminus C(p)$ is homeomorphic to an open m -ball if M is a m -manifold. This indicates that the topology of M is contained mostly in $C(p)$. We make this statement more precise using homology groups. For a topological space \mathbb{X} , let $H_j(\mathbb{X})$ denote the j -dimensional homology group defined over the field \mathbb{Z}_2 . The rank of $H_j(\mathbb{X})$ is called the j th Betti number of \mathbb{X} and denoted $\beta_j(\mathbb{X})$. In what follows we write $X_1 \approx X_2$ for two groups X_1 and X_2 if they are isomorphic. The following results relate topology of M to its cut locus [9, 21].

Proposition 2.1 *Let M be a compact Riemannian m -manifold without boundary and $p \in M$ be any point.*

1. $M \setminus \{p\}$ deformation retracts to $C(p)$ and $M \setminus C(p)$ deformation retracts to p ;
2. for $0 \leq j \leq m - 1$, $H_j(M) \approx H_j(C(p))$.

If the coefficient ring in the homology group is a field which is not necessarily \mathbb{Z}_2 , the second assertion remains true if M is orientable and is true only for $0 \leq j \leq m - 2$ if M is non-orientable.

3 Surface cut locus

We consider the case when $M \subset \mathbb{R}^k$ is a surface (2-manifold). The cut locus of an arbitrary smooth surface M could be structurally intractable in the sense that it may not even be triangulable. Fortunately, there is a large class called *real-analytic surfaces* that do not exhibit this pathological behavior. These surfaces can be described locally by a real analytic function (locally agrees with Taylor series expansion). It is known that the cut locus of any real-analytic, compact k -manifold is triangulable [4]. Henceforth, we assume that M is real-analytic.

3.1 Structural properties

A cut locus of a real-analytic surface is a *graph*, namely it consists of curve segments which are joined at vertices. Myers [22, 23] showed some more structural properties of cut loci. Let q be any point on a cut locus $C(p)$. In general, there could be one or more minimizing geodesics joining p and q . These geodesics may be separated or clumped together. To be precise consider a parameterization $\theta \mapsto \gamma_\theta$ where γ_θ is the geodesic γ with $\dot{\gamma}(0)$ making an angle θ with a fixed reference vector $v \in TM_p$. If q is a conjugate point to p , it is conceivable that there is an interval $[\theta_1, \theta_2]$ so that all geodesics γ_θ with $\theta \in [\theta_1, \theta_2]$ connect p and q . A remarkable result of Myers is that, this is not possible if M is a real-analytic surface unless the cut locus degenerates to a single point. We use this important structural property in our proofs. Actually, Myers [22, 23] proved more. Let the number of minimizing geodesics connecting p to a point $q \in C(p)$ be the *order* of q .

Proposition 3.1 *If $C(p)$ is not a single point, the order of a point $q \in C(p)$ is equal to the number of edges in $C(p)$ incident to q .*

Henceforth we assume that $C(p)$ is not a single point which happens only for geometric spheres and can be handled easily. One implication of Proposition 3.1 is that only finitely many minimizing geodesics connect a point to any point in its cut locus. Also, the degree of a vertex in the cut locus $C(p)$ is exactly equal to its order. In particular, *leaves*—the degree 1 vertices have exactly one geodesic coming into it. Generally, the cut locus $C(p)$, being a graph, contains cycles with tree structures attached to them. We call $q \in C(p)$ a *tree point* if either q is a leaf in $C(p)$ or $C(p) \setminus \{q\}$ contains a component whose relative closure in $C(p)$ is a tree. Otherwise, q is called a *cycle point*.

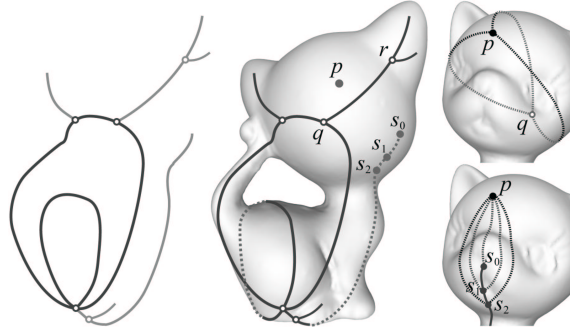


Figure 1: Cut locus on Kitten: cut locus drawn on a plane with tree points shaded lighter(left), cut locus embedded on the surface (right).

Notice that, for a tree point q which is not a leaf, $C(p) \setminus \{q\}$ contains a component contractible to q in M . In Figure 1, $C(p)$ has two cycles since the surface has genus 1. The points q, r, s_0, s_1, s_2 are tree points. Observe that even though q belongs to a cycle, it separates a tree and hence is a tree point by our definition. The order of q and r is three. The point s_0 is a conjugate point and its order is one. The two minimizing geodesics to s_2 are homotopic to each other. They separate a disk from the surface which contains all minimizing geodesics to the segment from s_0 to s_2 . However, for a cycle point this is not true. We show that two minimizing geodesics coming into a cycle point from p cannot be homotopic in M .

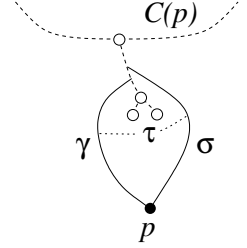
Proposition 3.2 *Suppose $q \in C(p)$ is a cycle point. A minimizing geodesic γ connecting p and q cannot be homotopic to any other minimizing geodesic $\sigma \neq \gamma$ connecting p and q .*

Proof. The two minimizing geodesics γ and σ meet only at p and q since two minimizing geodesics issued from p cannot meet in $M \setminus C(p)$. Therefore, γ and σ form a simple cycle in M . If γ and σ were homotopic, this cycle would bound an open disk, say D in M . If D does not intersect $C(p)$, we have an interval $[\theta_1, \theta_2]$ where $\gamma = \gamma_{\theta_1}$ and $\sigma = \gamma_{\theta_2}$ such that all minimizing geodesics $\{\gamma_\theta\}$, $\theta \in [\theta_1, \theta_2]$ connect p and q . This would violate proposition 3.1.

So, assume that $D \cap C(p)$ is non-empty. We claim that $D \cap C(p)$ cannot contain a cycle. For if there were a cycle in $C(p)$ embedded in a disk, it would contradict the fact that $M \setminus C(p)$ deformation retracts to p and hence is connected (Proposition 2.1). We are left with the only option of $D \cap C(p)$ being a non-empty tree. Therefore, q is a tree point by definition. But, this violates the assumption that q is a cycle point. It follows that γ and σ are not homotopic. \square

3.2 Geodesic spread

Recall that our goal is to compute the topology of M from a cut locus $C(p)$. Unfortunately, we cannot compute an approximation of the entire cut locus $C(p)$. Instead, we approximate a subset of it which still retains the topological information of $C(p)$. This subset can be defined in terms of a notion of *geodesic spread*, which we now develop. A subset of $C(p)$ retains the essential topology of $C(p)$ if it consists of points where two minimizing geodesics meet after spreading apart by an amount of at least $i(M)$, the injectivity radius of M . We formalize and prove this fact and then design an algorithm to approximate such a subset.



Spread *spd*. Let $\gamma : [0, t_0] \rightarrow M$ and $\sigma : [0, t_1] \rightarrow M$ be two minimizing geodesics parameterized by arc lengths which connect $p \in M$ to $\gamma(t_0) \in C(p)$ and $\sigma(t_1) \in C(p)$ respectively. Let $t_0 \leq t_1$. The distance $spd(\gamma, \sigma)$ is defined as

$$spd(\gamma, \sigma) = \max_{t \in [0, t_0]} \{d_M(\gamma(t), \sigma(t))\}.$$

This distance measures how far apart two geodesics get when traveling from p to the cut locus. In the figure above $spd(\gamma, \sigma) = \tau$.

Consider a function $\omega : C(p) \rightarrow \mathbb{R}$ where $\omega(q)$ is the maximum of $spd(\gamma, \sigma)$ over all pairs of minimizing geodesics γ, σ connecting p and q . For any $\tau \geq 0$, we also define $C_\tau(p) \subseteq C(p)$ as the set of points $\{q \in C(p)\}$ where $\omega(q) \geq \tau$. We aim to approximate a superset of $C_\tau(p)$ for some $\tau \leq i(M)$. The reason is that such a subset of $C(p)$ contains all information about the one dimensional homology group of M . To prove this fact, we establish first the following result.

Proposition 3.3 *Let γ_1, γ_2 be two minimizing geodesics connecting p and $q \in C(p)$. If $w(q) < i(M)$, γ_1 and γ_2 are homotopic.*

Proof. Consider the minimizing geodesic σ_t connecting $\gamma_1(t)$ and $\gamma_2(t)$, $0 \leq t \leq t_c$ where $\gamma_1(t_c) = \gamma_2(t_c) = q$, the cut point along γ_1 and γ_2 . We have $\sigma_t(0) = \gamma_1(t)$ and let $\sigma_t(b) = \gamma_2(t)$. Consider the map $f : \mathbb{R} \times [0, t_c] \rightarrow M$ given by $f(s, t) = \sigma_t(s)$. We show that the restriction of this map for $s \in [0, b]$ is smooth.

Let $S_p(M) \subset T_p(M) \times T_p(M)$ be the smooth 2-manifold defined by

$$S_p(M) = \{(v_1, v_2) \mid \|v_1\| = \|v_2\| = 1\}.$$

Let $\phi_p : \mathbb{R} \times S_p(M) \rightarrow M \times M$ be the smooth map defined by

$$(t, v_1, v_2) \mapsto (\exp_p(tv_1), \exp_p(tv_2)).$$

For two points $x, y \in M$ where $d_M(x, y) \leq i(M)$, let $u(x, y)$ denote the unit tangent vector of the minimizing geodesic between x and y at x . Now consider the map $\psi : [0, \infty) \times M \times M \rightarrow M$ restricted to the open set of $\{x, y\} \subseteq M \times M$ where $d_M(x, y) < i(M)$ and $\psi(s, x, y) = \exp_x(su(x, y))$. The map ψ is also smooth. Therefore, the composition $\psi \circ (id_{\mathbb{R}} \times \phi_p)$ is smooth ($id_{\mathbb{R}}$ is the identity on \mathbb{R}). Since $f(s, t) = \psi((id_{\mathbb{R}} \times \phi_p)(s, t, \dot{\gamma}_1(0), \dot{\gamma}_2(0)))$, one concludes that f is smooth.

Consider the continuous function $F : [0, 1] \times \mathbb{R} \rightarrow M$ given by $F(w, t) = f(wd_M(\gamma_1(t), \gamma_2(t)), t)$. We have $F(0, t) = \gamma_1(t)$ and $F(1, t) = \gamma_2(t)$. Thus, F is a homotopy between γ_1 and γ_2 proving the claim. \square

Combining Proposition 3.2 and Proposition 3.3 we conclude:

Proposition 3.4 *A point $q \in C(p)$ is a tree point if $\omega(q) < i(M)$.*

Next we show that any closed set containing cycle points captures the topology of $C(p)$. The closure is needed to take care of points such as q in Figure 1. Let $\text{Cl } Y$ denote the relative closure of a set $Y \subseteq C(p)$.

Proposition 3.5 *Let $X \subseteq C(p)$ be any closed set containing all cycle points of $C(p)$. Then, $H_1(X) \approx H_1(C(p))$.*

Proof. Let $Y \subseteq X$ be the set of all cycle points. Since X is closed, $\text{Cl } Y \subseteq X$. First we show that $H_1(\text{Cl } Y) \approx H_1(C(p))$. If $\text{Cl } Y = Y$, there is no tree point in the closure of cycle points. It means $C(p) = \text{Cl } Y$ since otherwise the only possibility is that $C(p)$ is disconnected which contradicts the fact that M is connected. So, assume that $Y \subset \text{Cl } Y$ and let y be any point in $\text{Cl } Y \setminus Y$. Since Y contains all cycle points, y is a tree point. The tree rooted at y trivially contracts to y . Contracting all such trees for all points $y \in \text{Cl } Y \setminus Y$, we are left with a set, say $Y' \supseteq \text{Cl } Y$ so that

$$H_1(Y') \approx H_1(C(p))$$

since contracting trees to points does not alter homology. We claim that $Y' = \text{Cl } Y$. If not, consider the set $Y'' = Y' \setminus \text{Cl } Y$. The set Y'' is not connected to $\text{Cl } Y$ and hence $C(p)$ is not connected, an impossibility. It follows that

$$H_1(Y') \approx H_1(\text{Cl } Y) \approx H_1(C(p)).$$

Observe that adding any subset of $C(p)$ to $\text{Cl } Y$ does not add any cycle. If it did, $\beta_1(C(p))$ would be larger than $\beta_1(\text{Cl } Y)$. Therefore, for any $X \supseteq \text{Cl } Y$, $\beta_1(X) = \beta_1(\text{Cl } Y)$ which proves that $\beta_1(X) = \beta_1(C(p))$, or $H_1(X)$ is isomorphic to $H_1(C(p))$. \square

As a corollary of Proposition 2.1, Proposition 3.4, and Proposition 3.5 we obtain:

Theorem 3.1 *For any closed set X where $C_\tau(p) \subseteq X \subseteq C(p)$, $\tau \leq i(M)$, we have $H_1(X) \approx H_1(C(p)) \approx H_1(M)$.*

Proof. The complement of X , $C(p) \setminus X$, contains only tree points by Proposition 3.4. Therefore, X contains all cycle points. Propositions 2.1 and 3.5 imply the conclusion immediately. \square

In our algorithm we approximate a superset of $C_\tau(p)$ for some $\tau \leq i(M)$ to honor Theorem 3.1. The algorithm approximates geodesics by shortest paths in an appropriate graph G spanned by a given point set $P \subset M$. For a point $q \in P$ and a graph G with vertices in P , let

$$\begin{aligned} \gamma_{pq}^G &= \text{shortest path between } p \text{ and } q \text{ in } G \text{ and} \\ d_G(p, q) &= \text{Euclidean length of } \gamma_{pq}^G. \end{aligned}$$

When the fixed source p is understood for all geodesics and shortest paths, we write

$$\gamma_q = \gamma_{pq} \text{ and } \gamma_q^G = \gamma_{pq}^G.$$

The algorithm approximates the distance $\text{spd}(\gamma_q, \gamma_s)$ between the two geodesics emanating from p by computing a distance similar to spd between the shortest paths γ_q^G and γ_s^G .

If this approximate distance is larger than a threshold τ , it selects end vertices q and s if they are close such as the ones in Figure 2 (the two shortest paths are approximating the two geodesics shown with dotted curves). If τ is relatively small compared to $i(M)$, the algorithm at least approximates an appropriate subset X of $C(p)$ which satisfies Theorem 3.1. At the same time the algorithm should not compute points far away from $C(p)$ even though it captures all points of X . Observe that if $\text{spd}(\gamma_q, \gamma_s)$ is not small, the two geodesics cannot come close unless they are near p or near a cut point. We use this observation to restrict all output points near $C(p)$. We need to define an open set containing $C(p)$ to make our statement precise. For $\eta > 0$, let

$$W_\eta(p) = \{\gamma(t) | t > t_c - \eta \geq 0 \text{ where } \gamma(t_c) \in C(p)\}$$

which is a space that deformation retracts to $C(p)$ along the minimizing geodesics originating from p .

Proposition 3.6 *For any $\eta > 0$ and $\tau > 0$ there exists a number $\nu = \nu(\eta, \tau) \geq 0$ so that if $x \in M \setminus W_\eta(p)$, there is a ν -neighborhood U of x where for any two points $u, u' \in U$, $\text{spd}(\gamma_u, \gamma_{u'}) < \tau$. Same statement also holds for the distance $d_H^E(\gamma_u, \gamma_{u'})$.*

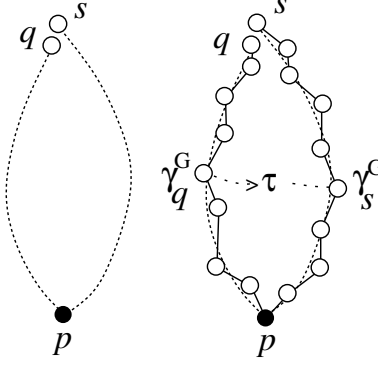


Figure 2: Geodesics approximated by shortest paths.

Proof. Let $C^0(\mathbb{R}, M)$ denote the space of all continuous functions from \mathbb{R} to M . Define $G : M \setminus W_\eta(p) \rightarrow C^0(\mathbb{R}, M)$ by taking $G(x) = \gamma_x$ where $\gamma_x(0) = p$. Since x is not in $C(p)$, γ_x is well defined. Let L be the length of γ_x . Define an open set

$$\mathcal{U}_\tau = \{\alpha : \mathbb{R} \rightarrow M \mid d_M(\alpha(t), \gamma_x(t)) < \tau \text{ for } 0 \leq t \leq L\}.$$

Since G is continuous by the continuity of geodesic flows, the inverse image $G^{-1}(\mathcal{U}_\tau)$ of the open set \mathcal{U}_τ is open. It follows that there is a ν -neighborhood of x contained in $G^{-1}(\mathcal{U}_\tau)$ as claimed. \square

3.3 Geodesic and spread approximation

We approximate the minimizing geodesics by shortest paths from p in an appropriate graph built on the point data that samples M . For this approximation to be good, we need that P sample M well.

Sampling condition. We say $P \subset M$ is an ε -sample if each point $x \in M$ has a point in P within geodesic distance of ε , that is, $d_M(x, P) \leq \varepsilon$.

For a point set $P \subseteq M$, let $G^\delta(P)$ denote the graph with vertices in P and edges that connect any two points $p, q \in P$ within Euclidean distance δ , that is, $d_E(p, q) \leq \delta$. Consider a sequence $\{P_n\}$ of point sets converging to M , that is, the sequence $\{\varepsilon_n\}$ approaches 0 where P_n is an ε_n -sample of M . For $\delta_n = \Theta(\sqrt{\varepsilon_n})$, define the sequence of graphs $\{G_n = G^{\delta_n}(P_n)\}$. For two points $p, q \in P_n$, let $\tilde{\gamma}^n = \gamma_{pq}^{G_n}$ be the shortest path between them in G_n . Let $W(p)$ be an open set containing the cut locus $C(p)$. We have the following claim.

Proposition 3.7 *Let $\{P_n\}$ be a sequence of ε_n -sample of M converging to it. For $\delta_n = \Theta(\sqrt{\varepsilon_n})$, let $\{G_n = G^{\delta_n}(P_n)\}$ be a sequence of graphs induced by $\{P_n\}$. For any open set $W(p)$ containing the cut locus $C(p)$, $p \in P_n$, the sequence of paths $\{\tilde{\gamma}^n\}$ between p and a point $q \in (M \setminus W(p)) \cap P_n$ converges uniformly to the unique minimizing geodesic γ between p and q in M .*

We prove the above proposition using a technique proposed by Hildebrandt, Polthier, and Wardetzky [19] to prove a similar result for convergence of geodesics on polyhedral surfaces. First, we need results on approximating geodesic distances. Actually, the proposition is proved by showing that the convergence in path lengths translates into a convergence in actual paths. Recall that $d_G(p, q)$ denote the length of the shortest path between two vertices p and q in a graph G .

Proposition 3.8 *Let p and q be two points as defined in Proposition 3.7. There exist two reals $\underline{\lambda}_n$ and $\overline{\lambda}_n$ so that*

$$(1 - \underline{\lambda}_n)d_M(p, q) \leq d_{G_n}(p, q) \leq (1 + \overline{\lambda}_n)d_M(p, q)$$

where $\overline{\lambda}_n, \underline{\lambda}_n \rightarrow 0$ as $n \rightarrow \infty$.

Proof. Let $p = v_0, v_1, \dots, v_k = q$ be the sequence of vertices on $\tilde{\gamma}^n$. Assuming ε_n sufficiently small, we have $\delta_n \leq \pi\rho(M)$ where $\rho(M)$ is the reach of M . Therefore,

$$d_E(v_i, v_{i+1}) \leq \delta_n \leq \pi\rho(M).$$

Under this condition we can apply Corollary 4 of [2] to claim

$$\begin{aligned} d_E(v_i, v_{i+1}) &\geq (1 - O(\delta_n))d_M(v_i, v_{i+1}) \\ &\geq (1 - O(\sqrt{\varepsilon_n}))d_M(v_i, v_{i+1}) \end{aligned}$$

which immediately gives

$$\begin{aligned} d_{G_n}(p, q) &\geq \sum_{i=0}^{k-1} (1 - O(\sqrt{\varepsilon_n}))d_M(v_i, v_{i+1}) \\ &\geq (1 - O(\sqrt{\varepsilon_n}))d_M(p, q). \end{aligned} \tag{1}$$

On the other hand, for $\varepsilon_n \leq \delta_n/4$, Theorem 2 of [2] provides

$$d_{G_n}(p, q) \leq (1 + 4\varepsilon_n/\delta_n)d_M(p, q). \tag{2}$$

In our case, since $\delta_n = \Theta(\sqrt{\varepsilon_n})$, the condition of $\varepsilon_n \leq \delta_n/4$ is satisfied for sufficiently small ε_n . We get

$$d_{G_n}(p, q) \leq (1 + O(\sqrt{\varepsilon_n}))d_M(p, q).$$

The proposition is established where both $\underline{\lambda}_n$ and $\overline{\lambda}_n$ are $O(\sqrt{\varepsilon_n})$ which goes to zero as ε_n goes to zero. \square

Proof. [Proposition 3.7] We will work on paths in M . To do so, we consider the path γ^n which consists of minimizing geodesics between all pairs of consecutive vertices on $\tilde{\gamma}^n$. Assuming that γ^n is arc length parameterized, we have

$$d_M(\gamma^n(t), \gamma^n(t')) \leq |t - t'| \tag{3}$$

for any t, t' in the domain of γ^n . We deduce from inequalities 1 and 2

$$|t - t'| \leq 2d_{G_n}(p, q) \leq 4d_M(p, q) \leq 4\text{diam}(M). \tag{4}$$

where $\text{diam}(M)$ is the diameter of M .

Let $k = 4\text{diam}(M)$ and $C^0([0, k], M)$ denote the space of all continuous functions $c : [0, k] \rightarrow M$. Interpret any path

$$c : [0, b] \rightarrow M, \quad b \leq k$$

as an element of $C^0([0, k], M)$ by considering $\tilde{c} : [0, k] \rightarrow M$ where

$$\tilde{c}(t) = \begin{cases} c(t) & \text{for } 0 \leq t \leq b \\ c(b) & \text{for } b \leq t \leq k. \end{cases}$$

Observe that the family $\{\gamma^n\}$ belongs to $C^0([0, k], M)$ due to the inequality 4. It follows from inequality 3 that the family $\{\gamma^n\}$ is equicontinuous. Also, the inequality 4 implies that the sequence $\{\gamma^n\}$ is uniformly bounded. Therefore, Arzelà-Ascoli theorem from functional analysis applies to establish that the set of accumulation points of $\{\gamma^n\}$ is not empty in the compact-open topology on $C^0([0, k], M)$.

Let γ be an accumulation point of $\{\gamma^n\}$. For a path $c : [0, b] \rightarrow M$, $b \leq k$, let $\ell(c)$ denote its length which is the supremum over all partitions of $Z = \{t_0 = 0 \leq t_1 \leq \dots \leq t_m = b\}$, that is, $\ell(c) = \sup_Z \sum_{i=1}^m d_M(c(t_{i-1}), c(t_i))$. We have

$$d_{G_n}(p, q) \leq \ell(\gamma^n) \leq (1 + O(\sqrt{\varepsilon_n}))d_{G_n}(p, q)$$

for small ε_n (apply the bound in the inequality 1). Then, Proposition 3.8 implies that $\ell(\gamma^n) \rightarrow d_M(p, q)$. The length functional $\ell : C^0([0, k], M) \rightarrow [0, \infty]$ is lower semicontinuous. Therefore,

$$\ell(\gamma) \leq \liminf \ell(\gamma^n) = d_M(p, q).$$

Hence γ is a minimizing geodesic connecting p and q . Since q lies outside an open neighborhood of the cut locus of p , there is a unique such geodesic between p and q meaning that the sequence $\{\gamma^n\}$ converges to the minimum geodesic γ between p and q . Moreover, Arzelà-Ascoli theorem says that this convergence is uniform. One can deduce from Proposition 3.8 that $d_{G_n}(p, q)$, the length of $\tilde{\gamma}^n$ approaches $\ell(\gamma^n)$ as $n \rightarrow \infty$. It follows that $\tilde{\gamma}^n$ converges to γ uniformly as well. \square

Corollary 3.1 For any $\mu > 0, \eta > 0$, there exists $\varepsilon = \varepsilon(\mu, \eta) > 0$ so that if P is an ε -sample and p and q are two points in P with $q \notin W_\eta(p)$, the shortest path γ_q^G between p and q in $G^{\Theta(\sqrt{\varepsilon})}(P)$ satisfies $d_H^E(\gamma_q^G, \gamma_q) \leq \mu$ where γ_q is the minimizing geodesic connecting p and q in M .

The above corollary relates shortest paths between vertices to the minimizing geodesics between them. We can state a similar fact for all minimizing geodesics issued from p .

Proposition 3.9 For any $\mu > 0, \eta > 0$, there exists $\varepsilon > 0$ so that if P is an ε -sample, then for any minimizing geodesic γ_q issued from p with $q \notin W_\eta(p)$ there is a shortest path γ_q^G originating from p in $G^{\Theta(\sqrt{\varepsilon})}(P)$ where $d_H^E(\gamma_q, \gamma_q^G) \leq \mu$.

Proof. Let $\bar{q} \in (M \setminus W_\eta(p)) \cap P$ be the closest point to q in terms of geodesic distance. For any μ and η , we can choose ε to be small enough so that $d_H^E(\gamma_q, \gamma_{\bar{q}}) < \mu/2$ by appealing to Proposition 3.6. The claim follows since $d_H^E(\gamma_{\bar{q}}, \gamma_{\bar{q}}^G) \leq \mu/2$ can be assumed for sufficiently small ε . \square

Now we show how we approximate the spread. Once we approximate the minimizing geodesics with the shortest paths in $G = G^\delta(P)$, $\delta = \Theta(\sqrt{\varepsilon})$, we can approximate the spread $\text{spd}(\gamma_q, \gamma_s)$ between two minimizing geodesics to q and s respectively. For this we consider the shortest paths γ_q^G and γ_s^G , and for each $v \in \gamma_q^G$ we find the set of vertices $V \subset \gamma_s^G$ that have nearly equal distance from the root p as v . The shortest paths from v to all vertices of γ_s^G in V approximate the distance $d_M(\gamma_q(t), \gamma_s(t))$ where $v = \gamma_q(t)$. We take the largest one among these paths to approximate the length $d_M(\gamma_q(t), \gamma_s(t))$. Let the length of this largest path be $\ell(v)$. Define

$$\text{spd}_G(\gamma_q^G, \gamma_s^G) = \max_{v \in \gamma_q^G} \{\ell(v)\}$$

which is computed by APPROXSPD. See Figure 3 for an illustration of the computed spread values.

APPROXSPD(γ_q^G, γ_s^G)

1. Assume $G = G^\delta(P)$ is available; $\ell_{max} := 0$;
2. for each vertex v on γ_q^G do
 - (a) Determine the vertex set $V = \{w\}$ so that $w \in \gamma_s^G$ and $d_G(p, v) - \delta \leq d_G(p, w) \leq d_G(p, v) + \delta$;
 - (b) Compute the largest length ℓ of the shortest paths from v to $w \in V$ in $G^\delta(P)$;
 - (c) if $(\ell > \ell_{max})$ $\ell_{max} := \ell$;
3. Return ℓ_{max} .

Proposition 3.10 For any $\phi > 0$ and $\eta > 0$ there is an $\varepsilon > 0$ so that if P is an ε -sample and $\delta = \Theta(\sqrt{\varepsilon})$ then for $q, s \in M \setminus W_{2\eta}(p)$, APPROXSPD(γ_q^G, γ_s^G) returns $\text{spd}_G(\gamma_q^G, \gamma_s^G)$ where $\text{spd}(\gamma_q, \gamma_s) - \phi < \text{spd}_G(\gamma_q^G, \gamma_s^G) < \text{spd}(\gamma_q, \gamma_s) + \phi$.

Proof. Let v and v' be the vertices on γ_q^G and γ_s^G respectively realizing $\text{spd}_G(\gamma_q^G, \gamma_s^G)$, that is, $d_G(v, v') = \text{spd}_G(\gamma_q^G, \gamma_s^G)$. For convenience, we write $a \in b \pm c$ if $b - c \leq a \leq b + c$. Notice that $a \in b \pm c$ implies $b \in a \pm c$. If ε is chosen sufficiently small, all vertices on γ_q^G and γ_s^G can be assumed to be outside $W_\eta(p)$ since q and s are outside $W_{2\eta}(p)$. This means we can apply Corollary 3.1.

Let \bar{v} be the closest point on γ_q to v . We have seen $d_G(p, v) \in d_M(p, v) \pm O(\sqrt{\varepsilon})d_M(p, v)$. Since $d_E(v, \bar{v}) \leq \mu$ (Corollary 3.1), we have

$$d_G(p, v) \in d_M(p, \bar{v}) \pm O(\sqrt{\varepsilon})d_M(p, v) + \mu. \quad (5)$$

We also have $d_M(p, v') \in d_G(p, v') \pm O(\sqrt{\varepsilon})d_G(p, v')$. The algorithm ensures $d_G(p, v') \in d_G(p, v) \pm \delta$ from which we get $d_M(p, v') \in d_G(p, v) \pm O(\sqrt{\varepsilon})d_G(p, v') + \delta$. Combining previous observations with the inequality 5, we get

$$\begin{aligned} d_M(p, v') \in d_M(p, \bar{v}) &\pm O(\sqrt{\varepsilon})d_M(p, v) \\ &+ \sqrt{\varepsilon}d_G(p, v') + \mu + \delta \end{aligned} \quad (6)$$

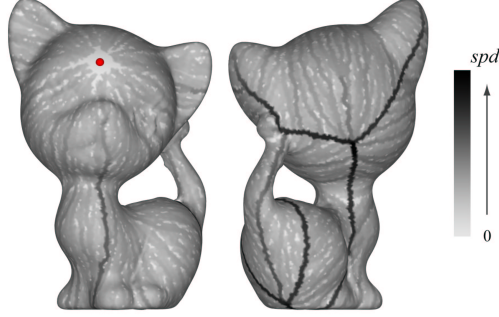


Figure 3: Pairs of vertices connected in G^δ are shaded according to the spread value computed for the two shortest paths joining the source to the vertices in the pairs; the darker the shades, the larger is the spread. Notice that pairs near the cut locus have distinctly large spread values.

Let \bar{v}' denote the closest point on γ_s to v' and let $z = \gamma_s(t)$ if $\bar{v} = \gamma_q(t)$. Using $d_E(v', \bar{v}') \leq \mu$ (Corollary 3.1) we obtain from the inequality 6

$$\begin{aligned} d_M(p, \bar{v}') &\in d_M(p, \bar{v}) \pm O(\sqrt{\varepsilon}d_M(p, v) \\ &\quad + \sqrt{\varepsilon}d_G(p, v') + \mu + \delta) \\ &= d_M(p, z) \pm O(\sqrt{\varepsilon}d_M(p, v) \\ &\quad + \sqrt{\varepsilon}d_G(p, v') + \mu + \delta). \end{aligned} \tag{7}$$

Since both \bar{v}' and z belong to γ_s , the inequality 7 provides

$$d_M(z, \bar{v}') = O(\sqrt{\varepsilon}d_M(p, v) + \sqrt{\varepsilon}d_G(p, v') + \mu + \delta)$$

from which we obtain

$$d_M(z, v') = O(\sqrt{\varepsilon}d_M(p, v) + \sqrt{\varepsilon}d_G(p, v') + \mu + \delta).$$

It follows

$$\begin{aligned} d_M(\bar{v}, z) &\in d_M(v, v') \pm (d_M(v', z) + d_M(v, \bar{v})) \\ &= d_M(v, v') \pm O(\sqrt{\varepsilon}d_M(p, v) \\ &\quad + \sqrt{\varepsilon}d_G(p, v') + \mu + \delta) \\ &= d_G(v, v') \pm O(\sqrt{\varepsilon}d_G(v, v') + \sqrt{\varepsilon}d_M(p, v) \\ &\quad + \sqrt{\varepsilon}d_G(p, v') + \mu + \delta). \end{aligned}$$

It follows that there is a ϕ_1 which goes to zero as ε does where $d_G(v, v') \leq d_M(\bar{v}, z) + \phi_1 \leq \text{spd}(\gamma_q, \gamma_s) + \phi_1$.

A very similar proof can show that there are two vertices $u \in \gamma_q^\varepsilon, u' \in \gamma_s^\varepsilon$ with $d_G(p, u) \in d_G(p, u') \pm \delta$ so that $\text{spd}(\gamma_q, \gamma_s) \leq d_G(u, u') + \phi_2$. Since $d_G(u, u') \leq d_G(v, v')$ by the definition of v and v' , we have $\text{spd}(\gamma_q, \gamma_s) \leq d_G(v, v') + \phi_2$ where $\phi_2 \rightarrow 0$ as $\varepsilon \rightarrow 0$. For any given ϕ , one can satisfy $\phi > \max\{\phi_1, \phi_2\}$ by choosing ε sufficiently small. Then, we obtain the result as claimed. \square

4 Cut locus approximation

The algorithm CUTLOCUS below implements the following strategy. It selects all pair of points $q, s \in P$ which are close and admit shortest paths γ_q^ε and γ_s^ε respectively in $G = G^\delta(P)$ where $\text{spd}_G(\gamma_q^\varepsilon, \gamma_s^\varepsilon)$ is more than a threshold.

CUTLOCUS($p \in P, \kappa, \xi, \delta$)

1. Compute the graph $G = G^\delta(P)$; $C := \emptyset$;
2. Compute shortest paths from a point $p \in P$ to all vertices in $G^\delta(P)$;
3. for each $q \in P - \{p\}$ do
 for each pair $q, s \in P$ with $d_E(q, s) \leq \xi$ do
 if APPROXSPD(γ_q^G, γ_s^G) $\geq \kappa$
 then $C := C \cup \{q, s\}$;
4. Return C .

4.1 Justification

Proposition 4.1 *For a sufficiently small positive η and any positives μ and ϕ , there is an $\varepsilon > 0$ so that if P is ε -sample, $\delta = \Theta(\sqrt{\varepsilon})$ and $\xi > 2(\mu + \delta + \eta)$, then the following is true: CUTLOCUS($p \in P, \kappa, \xi, \delta$) computes a sample point $q \in P$ for each $x \in C_\tau(p)$ where $d_E(q, x) \leq \mu + \eta + \delta$ and $\tau \geq \kappa + \phi + O(\mu)$.*

Proof. Let $x \in C_\tau(p)$ for some $\tau > 0$. There are two geodesics γ_1 and γ_2 connecting p and x where $\text{spd}(\gamma_1, \gamma_2) \geq \tau$. Let x_1 and x_2 be two points in γ_1 and γ_2 respectively on the boundary of $W_\eta(p)$. By Proposition 3.9, there is an $\varepsilon > 0$ so that there are two shortest paths γ_1^G and γ_2^G in the graph $G^\delta(P)$ for $\delta = \Theta(\sqrt{\varepsilon})$ where $d_H^E(\gamma_i^G, \gamma_{x_i}) = O(\mu)$ for any η and μ .

Let v_1 and v_2 be the closest vertices to x_1 and x_2 respectively on the paths γ_1^G and γ_2^G . We get

$$d_E(v_1, x) \leq d_E(v_1, x_1) + d_E(x_1, x) \leq \mu + \delta + \eta.$$

It follows that $d_E(v_1, v_2) \leq 2(\mu + \delta + \eta)$. Assuming $\xi > 2(\mu + \delta + \eta)$, the shortest paths $\gamma_{v_1}^G$ and $\gamma_{v_2}^G$ are checked for their distance by APPROXSPD.

We have $\text{spd}(\gamma_{v_1}, \gamma_{v_2}) \geq \text{spd}(\gamma_{x_1}, \gamma_{x_2}) - O(\mu)$ since both γ_{v_i} and γ_{x_i} makes $O(\mu)$ distance with $\gamma_{v_i}^G$. For a sufficiently small $\eta > 0$, $\text{spd}(\gamma_{x_1}, \gamma_{x_2}) = \text{spd}(\gamma_1, \gamma_2) \geq \tau$ by definition of x being in $C_\tau(p)$. Thus, $\text{spd}(\gamma_{v_1}, \gamma_{v_2}) \geq \tau - O(\mu)$. APPROXSPD($\gamma_{v_1}^G, \gamma_{v_2}^G$) returns a value more than $\tau - O(\mu) - \phi$ by Proposition 3.10. The vertices v_1 and v_2 are output by CUTLOCUS if $\tau - O(\mu) - \phi \geq \kappa$, or equivalently if $\tau \geq O(\mu) + \phi + \kappa$. Either of v_1 and v_2 can be taken as q since both of them satisfy the stated properties of q in the lemma. \square

Proposition 4.2 *For any $\eta > 0$, there exist $\varepsilon > 0$ and $\xi > 0$ so that the following holds. Let v be a vertex computed by CUTLOCUS($p \in P, \kappa, \xi, \delta$) where P is an ε -sample of M and $\delta = \Theta(\sqrt{\varepsilon})$. There is a point $x \in C(p)$ so that $d_E(x, v) \leq 2\eta$.*

Proof. If $d_E(v, C(p)) \leq 2\eta$, we are done. So, assume otherwise, that is, v lies outside the 2η -neighborhood of $C(p)$ in \mathbb{R}^k . Since v is computed by CUTLOCUS, there is another vertex v' computed by CUTLOCUS so that $d_E(v, v') \leq \xi$ and the shortest paths γ_v^G and $\gamma_{v'}^G$ satisfy $\text{spd}_G(\gamma_v^G, \gamma_{v'}^G) \geq \kappa$. By Proposition 3.10,

$$\text{spd}(\gamma_v, \gamma_{v'}) > \text{spd}_G(\gamma_v^G, \gamma_{v'}^G) - \phi$$

for any ϕ as long as ε is sufficiently small. Therefore, we have $\text{spd}(\gamma_v, \gamma_{v'}) > \kappa - \phi$.

Observe that we can assume $d_E(v', C(p)) > \eta$ since otherwise v' satisfies the lemma. We want to apply Proposition 3.6 to γ_v and $\gamma_{v'}$ with $\tau = \kappa - \phi$. For this τ , let $\nu = \nu(\tau, \eta)$ satisfy the proposition. If we choose $\xi \leq \nu$, one should have $\text{spd}(\gamma_v, \gamma_{v'}) < \tau$ reaching a contradiction. \square

Theorem 4.1 *For any $\varepsilon' > 0$, there is an $\varepsilon > 0$ and a $\xi > 2\varepsilon'$ so that if P is an ε -sample, $\delta = \Theta(\sqrt{\varepsilon})$, and $0 < \kappa < i(M) - O(\varepsilon')$, then $P' \subset P$ returned by CUTLOCUS(p, κ, ξ, δ) has $d_H^E(P', X) = O(\varepsilon')$ where $C_\tau \subseteq X \subseteq C(p)$ for $\tau < i(M)$.*

Proof. Observe that for any $\varepsilon' > 0$ we can satisfy $\varepsilon' > \mu + \delta + \eta + \phi$ by assuming ε to be sufficiently small. Then, if ξ is chosen where $\xi > 2\varepsilon'$, we have $\xi > 2(\mu + \delta + \eta)$. Therefore, we can apply Proposition 4.1 to claim that for each point $x \in C_\tau(p)$, $\text{CUTLOCUS}(p, \kappa, \xi, \delta)$ outputs a point with distance $d_E(x, p) = O(\varepsilon')$ where $\tau = \kappa + \phi + O(\mu)$. When κ falls into the stated range, we have $\tau < i(M)$.

From Proposition 4.2 we get that if ξ is sufficiently small, all points computed by CUTLOCUS are within $2\eta = O(\varepsilon')$ distance of $C(p)$. If ε is small enough, such a ξ can be chosen satisfying all constraints. Combining all, we get that CUTLOCUS computes a set P' where $d_H^E(P', X) = O(\varepsilon')$ and $C_\tau(p) \subseteq X \subseteq C(p)$ for $\tau < i(M)$. \square

5 Computing homology

Our algorithm for homology computation first estimates a sample density parameter with which we build a graph where shortest paths are computed. Rips complexes are constructed with an input parameter on the point set approximating a cut locus. Persistent Betti numbers are computed on these Rips complexes.

Estimating density. Observe that we need an estimate of ε to construct the graph $G^\delta(P)$ since we set $\delta = \Theta(\sqrt{\varepsilon})$. We estimate ε by using a procedure that was suggested in [3] to build a sequence of subsamples in the context of computing witness complexes [25].

Let $\{p_0\} = L_0 \subset L_1 \subset \dots \subset L_k = P$ where $L_{i+1} = L_i \cup \{p_{i+1}\}$ with p_{i+1} being the furthest point in $P \setminus L_i$ from L_i , that is, $p_{i+1} = \operatorname{argmax}_{q \in P \setminus L_i} d_E(q, L_i)$. Define $\tilde{\varepsilon}_i = d_E(p_{i+1}, L_i)$. We show that $\tilde{\varepsilon}_i$ approximates a sampling density ε_i defined as follows. A sample $L_i \subset M$ is a *tight ε_i -sample* if L_i is an ε_i -sample of M and there is an $x \in M$ for which $d_M(x, L_i) = \varepsilon_i$.

To prove that $\tilde{\varepsilon}_i$ approximates ε_i we need Proposition 5.2 which in turn uses Lemma 3 of [2]. Let $\frac{1}{r_0} = \max_{\gamma, t} \{\|\dot{\gamma}(t)\|\}$ where γ varies over all unit speed geodesics in M and $t \in \mathbb{R}$.

Proposition 5.1 ([2]) *For any two points $p, q \in M$, if $d_M(p, q) \leq \pi r_0$, then $d_E(p, q) \geq 2r_0 \sin(\frac{d_M(p, q)}{2r_0})$.*

Recall that the reach $\rho(M)$ is the smallest distance between M and its medial axis. The reach $\rho(M)$ bounds r_0 from below. See [13] for a proof of this fact.

Proposition 5.2 *For any two points $p, q \in M$, if $d_M(p, q) \leq \rho(M)/2$, then $d_E(p, q) \geq \frac{9}{10} d_M(p, q)$.*

Proof. First, observe that $d_M(p, q) \leq \frac{\rho(M)}{2} \leq \pi r_0$ which allows us to apply Proposition 5.1. Second, $\sin(t) \geq t - t^3/6$ for $t \geq 0$. Plugging this into the bound given by Proposition 5.1 and writing $\ell = d_M(p, q)$, we get

$$d_E(p, q) \geq (1 - \frac{\ell^2}{24r_0^2})\ell \geq (1 - \frac{\ell^2}{24\rho(M)^2})\ell.$$

Since $\ell \leq \rho(M)/2$, we have

$$d_E(p, q) \geq (1 - \frac{1}{96})\ell \geq \frac{9}{10} d_M(p, q).$$

\square

Notice that the choice of the factor $\frac{9}{10}$ is a little arbitrary. We could have taken the factor $\frac{95}{96}$ which would tighten other constants slightly.

Proposition 5.3 *Let $\{L_i\}$ be the sequence of subsamples as described above. If $L_i \subset P$ is a tight ε_i -sample and P is an ε -sample of M respectively, then for $\varepsilon < \varepsilon_i \leq \rho(M)/2$, one has $\frac{9}{10}(\varepsilon_i - \varepsilon) \leq \tilde{\varepsilon}_i \leq \varepsilon_i$.*

Proof. Consider a point $x \in M$ so that $d_M(x, L_i) = \varepsilon_i$. Since L_i is a tight ε_i -sample such a point exists. Let w be the closest point to x in $P \setminus L_i$. We claim that w is also the closest point to x in P . If not, there is a point in L_i which is closest to x in P . Then, $d_M(x, L_i) \leq \varepsilon$ contradicting that $\varepsilon < \varepsilon_i$.

We have $d_M(w, x) \leq \varepsilon$ since P is an ε -sample. Let p be the closest point to w in L_i . Then,

$$d_M(w, p) \geq d_M(x, p) - \varepsilon \geq d_M(x, L_i) - \varepsilon = \varepsilon_i - \varepsilon.$$

Since L_i is an ε_i -sample, $d_M(w, p) \leq \varepsilon_i \leq \frac{\rho(M)}{2}$. We can apply Proposition 5.2 to claim $d_E(w, p) \geq \frac{9}{10}d_M(w, p)$. Then, we have $\tilde{\varepsilon}_i \geq d_E(w, p) \geq \frac{9d_M(w, p)}{10} \geq \frac{9}{10}(\varepsilon_i - \varepsilon)$. This proves the lower bound on $\tilde{\varepsilon}_i$.

To prove the upper bound, consider the pair (u, p) , $p \in L_i$, $u \in P \setminus L_i$ which realizes the distance $\tilde{\varepsilon}_i$. Since L_i is an ε_i -sample and $u \in M \setminus L_i$, $\tilde{\varepsilon}_i \leq d_M(u, p) \leq \varepsilon_i$. \square

Now we have all ingredients to compute $\beta_1(M)$ from a dense sample P of M . For any compact set $\mathbb{X} \subseteq \mathbb{R}^k$, let \mathbb{X}^α denote its α -offset in \mathbb{R}^k , that is,

$$\mathbb{X}^\alpha = \{x \in \mathbb{R}^k \mid \inf_{y \in \mathbb{X}} d_E(x, y) \leq \alpha\}.$$

Let C denote the set of critical points of the distance function d_E restricted to the domain of \mathbb{X} . The distance

$$\text{wfs}(\mathbb{X}) = \inf_{x \in \mathbb{X}} \inf_{c \in C} d_E(x, c)$$

is called the weak feature size of \mathbb{X} [6].

Rips complexes. The α -Rips complex of a point set $P \subset \mathbb{R}^k$ is defined as a simplicial complex $\mathcal{R}^\alpha(P)$ where a simplex σ with vertices in P is in $\mathcal{R}^\alpha(P)$ if and only if all edges of σ has length at most 2α . Chazal and Oudot [7] show that the j th Betti number $\beta_j(\mathbb{X}) = \text{rank } H_j(\mathbb{X})$ can be computed from Rips complexes as follows.

Let $P \subseteq \mathbb{X}$ be a point sample of \mathbb{X} with the Hausdorff distance $d_H^E(P, \mathbb{X}) \leq \varepsilon$. The natural inclusion $\mathcal{R}^\alpha(P) \xhookrightarrow{\iota} \mathcal{R}^{4\alpha}(P)$ also induces a homomorphism ι^* at the homology level

$$H_j(\mathcal{R}^\alpha(P)) \xrightarrow{\iota^*} H_j(\mathcal{R}^{4\alpha}(P)).$$

The integer

$$\beta_j^{\alpha, 4\alpha}(P) = \text{rank}(\text{image } \iota^*)$$

denotes the j th persistent Betti number given by the inclusion of $\mathcal{R}^\alpha(P)$ into $\mathcal{R}^{4\alpha}(P)$. Persistent Betti numbers can be computed by the persistence algorithm pioneered by Edelsbrunner, Zomorodian and Letscher [15] and extended later by Zomorodian and Carlsson [27]. It is shown in [7] that for any offset \mathbb{X}^λ , $0 < \lambda \leq \text{wfs}(\mathbb{X})$, one has $\beta_j(\mathbb{X}^\lambda) = \beta_j^{\alpha, 4\alpha}(P)$ if $2\varepsilon \leq \alpha \leq \frac{1}{4}(\text{wfs}(\mathbb{X}) - \varepsilon)$. After computing the complexes $\mathcal{R}^\alpha(P)$ and $\mathcal{R}^{4\alpha}(P)$ one can compute the persistent Betti numbers by following the persistence algorithm [15, 27] on a filtration that adds the simplices of $\mathcal{R}^{4\alpha}(P) \setminus \mathcal{R}^\alpha(P)$ to $\mathcal{R}^\alpha(P)$.

Algorithm. We want to follow the same approach for the cut locus that we approximate by CUTLOCUS. Let $X \subseteq C(p)$ be the closed set approximated by the point set $L'_i \subseteq L_i$ that CUTLOCUS computes. By Theorem 3.1 and Theorem 4.1, $H_1(X) \approx H_1(M)$ if appropriate parameter values are passed to CUTLOCUS and L_i is sufficiently dense. Let

$$\rho_1 = \sup_{\alpha} (H_1(X^\alpha) \approx H_1(X)).$$

Following [11], one may call ρ_1 the *first homological feature size* of X . It turns out that $\text{wfs}(X)$ is bounded above by ρ_1 . The technique of [7] can be used to show that $\beta_1^{\alpha, 4\alpha}(L'_i)$ is equal to $\beta_1(X)$ if $4\varepsilon'_i \leq \alpha \leq \frac{1}{4}(\rho_1 - \varepsilon'_i)$ where $d_H^E(L'_i, X) \leq \varepsilon'_i$.

For a large range of i , L_i remains a dense sample of M . Also, by Proposition 5.3 the estimated sampling density $\tilde{\varepsilon}_i$ of L_i follows closely its actual density ε_i for a large range of i as long as ε_i remains larger than ε . Therefore, to estimate ε_i properly, we consider all L_i s iteratively in the algorithm. Theorem 4.1 requires that $d_H^E(L'_i, X)$ is at most $\frac{\xi}{2}$. If ξ satisfies $3\xi \leq \frac{1}{4}(\rho_1 - \varepsilon'_i)$, we have that $4\varepsilon'_i \leq 3\xi \leq \frac{1}{4}(\rho_1 - \varepsilon'_i)$. The constraints on ξ can be satisfied if ε_i is small enough. Then, taking $\alpha = 3\xi$ we can compute $\beta_1^{3\xi, 12\xi}(L'_i)$ which equals $\beta_1(X)$.

TOPODATA(P, κ, ξ)

1. Initialize $L = \emptyset$;
2. while $L \neq P$ do
 - (a) compute $p := \operatorname{argmax}_{q \in P} d_E(q, L)$;
 $L := L \cup \{p\}$; $P := P \setminus \{p\}$;
 - (b) compute $\tilde{\varepsilon} := d_E(p, L)$;
 - (c) Compute $L' := \text{CUTLOCUS}(p \in L, \kappa, \xi, \sqrt{\tilde{\varepsilon}})$;
 - (d) Output $\beta_1^{3\xi, 12\xi}(L')$ by considering the inclusion $\mathcal{R}^{3\xi}(L') \hookrightarrow \mathcal{R}^{12\xi}(L')$;

The output of TOPODATA can be plotted against the filtration of the point sample and the most persistent Betti numbers can be selected as outlined in [7]. Since $\tilde{\varepsilon}_i$ estimates ε_i for a large range of i , the persistent Betti numbers computed in step 2(d) remain stable for a large interval assuming that the original sample is dense enough for L'_i to remain dense for X for a large range of i . Notice that we do not have any relation established between ρ_1 , the first homological feature size of X and the injectivity radius of M . It is conceivable that a point sample which is dense for M may not provide dense enough sample for X . For this we assume that the original sample is so dense that L'_i remains dense for X for a large range of i . In particular, we assume that ε is sufficiently small so that ξ can be chosen sufficiently small satisfying $2\varepsilon'_i < \xi \leq \frac{1}{12}(\rho_1 - \varepsilon'_i)$ for Theorem 4.1 to hold and for computation of $\beta_1(X)$ to remain correct.

Time complexity. Let P have n sample points. First, we determine the time complexity of the algorithm CUTLOCUS. Computation of the graph G^δ cannot take more than $O(n^2)$ time since it involves checking pairwise distances of points in P . Computation of the shortest paths from a source in step 2 of CUTLOCUS takes $O(n^2)$ time. For step 3 we need to determine shortest paths between different vertices in G^δ . We compute all pairs shortest paths in G^δ and keep the pairwise distances in a matrix form. Once this is computed, step 3 of CUTLOCUS can be implemented in $O(n^3)$ time. Therefore, CUTLOCUS runs in $O(n^3)$ time.

In TOPODATA steps 2(a-b) can be performed in $O(n^2)$ time with a straightforward pairwise distance computations. Step 2(c) takes $O(n^3)$ time as we argued. Since persistence algorithm takes time cubic in its input size, step 2(d) takes $O(k^3)$ time where k is the number of simplices in $\mathcal{R}^{12\xi}(L')$. Since $C(p)$ is smooth everywhere except at its vertices, the analysis of [7] can be carried out to claim that $k = O(n)$. It implies that the step 2(d) takes $O(n^3)$ time. Accounting for all iterations, we obtain that TOPODATA runs in $O(n^4)$ time.

In theory we do not gain any advantage by computing an approximation to the cut locus since in the worst case a sample may have most points concentrated near a cut locus that is being approximated. However, this is too pathological to happen in practice. In fact, for a uniform distribution, a cut locus of length ℓ has roughly $\frac{\ell}{\varepsilon}$ points whereas the surface with area A has roughly $\frac{A}{\varepsilon^2}$ points implying a reduction by a factor of $\frac{\ell}{\varepsilon}$. Our experiments in 3D shows that the number of points are drastically reduced by cut locus approximation, see Table 1.

6 Experiments and conclusions

We implemented the algorithms CUTLOCUS and TOPODATA and ran them on an Intel Xeon 2.66GHz, 4GB RAM machine. We deviated from theory slightly in the implementation. If we let $\xi = 2\tilde{\varepsilon}$, it satisfies the required constraints if ε is sufficiently small. Instead, in the implementation we take $\xi = \tilde{\varepsilon}$ to reduce the sizes of the Rips complexes. Also, we take κ a multiple of ξ (see Table 1). Instead of computing $\beta^{3\xi, 12\xi}(\cdot)$ we compute $\beta^{\xi, 2\xi}(\cdot)$ which gives correct result in all cases that we tested.

Figure 4 shows some of the results. We also provide the time data for different steps of the algorithms in Table 1. We observe that the point set output by CUTLOCUS is much smaller than the input point set. As a result the sizes of the Rips complexes become much smaller as the Table 1 shows. Consequently building the Rips complexes and computing the persistent Betti numbers from them take less time. The gain in time outweighs the extra time required to compute an approximation to the cut locus. In Table 1 we also show the times for the case when $L = P$.

Model	#v	Locus	κ	L-Rip(#e,#f)	L-Rip	L-Perst	G-Rip(#e,#f)	G-Rip	G-Perst
Torus	9.4k	1.08	8 ξ	3.9k, 13.3k	0.01	0.04	219.6k, 1897.7k	9.63	44.85
Kitten	31.3k	12.05	6 ξ	13.2k, 54.4k	0.07	0.18	940.5k, 10569.9k	81.11	507.15
2-Torus	39.7k	24.33	8 ξ	14.1k, 47.7k	0.05	0.15	873.3k, 7063.5k	32.97	663.77
Genus3	63.9k	81.26	8 ξ	19.2k, 74.3k	0.08	0.24	1709.7k, 17018.2k	112.05	1496.45
Botijo	68.4k	78.46	8 ξ	32.5k, 138.2k	0.19	0.48	2016.9k, 22232.8k	173.80	1004.03
Mother	75.1k	119.61	8 ξ	25.9k, 92.3k	0.11	0.31	1715.5k, 14437.1k	74.27	2712.46
Hip	104.2k	492.07	3 ξ	103.4k, 439.4k	0.74	1.62	2202.7k, 17231.1k	203.88	4727.57
Pegasus	141.5k	519.84	4 ξ	95.9k, 443.1k	0.73	1.41	2979.8k, 28279.1k	3900.71	7728.63

Table 1: #v column denotes the number of vertices for each point cloud. All times are in seconds. Locus column denotes the time for computing the cut locus with parameter passed to TOPODATA shown in κ column. The two numbers in L-Rip(#e,#f) column are the number of edges and faces of the Rips complex of the subsample approximating the cut locus. L-Rip and L-Perst column denote the time for computing this Rips complex and time for running persistence on this Rips complex. G-Rip(#e,#f), G-Rip and G-Perst columns have the same meaning as the previous three columns but correspond to the entire point cloud instead of the points approximating the cut locus.

A natural extension of our work would be to apply the approach to data sampled from high dimensional manifolds. Our algorithm applies to these cases straightforwardly since it only involves computing distances on shortest path graphs. However, we do not have a proof of correctness at the moment. We require a generalization of Theorem 3.1. This needs a generalization of the definitions of tree and cycle points. Also, we would like to prove a stronger version of Theorem 4.1 where the Hausdorff distance bound is in terms of ε instead of ε' . This would require strengthening of Propositions 3.6 and 3.7 in terms of ε .

Notice that the persistence algorithm provides generators for homology classes in addition to their ranks. Therefore, one may compute a set of cycles from the input point cloud data that represent a basis of the first homology group of the sampled surface. However, these cycles are not guaranteed to be optimal or close to optimal in terms of lengths. We plan to address this issue in future work.

Acknowledgments

We thank the anonymous referees and Frédéric Chazal whose comments helped greatly improve this paper. We also thank Professor Dan Burghlea from the OSU Mathematics Department who explained some of the concepts from differential geometry used in this paper. This work is supported by NSF funded research grant CCF-0830467.

References

- [1] M. Belkin, P. Niyogi. Towards a Theoretical Foundation for Laplacian Based Manifold Methods. *J. Comput. System Sci.*, to appear.
- [2] M. Bernstein, V. de Silva, J. Langford, and J. Tenenbaum. Graph approximations to geodesics on embedded manifolds. *Tech Report*, Dept. Psychology, Stanford University, USA, 2000. Available at <http://isomap.stanford.edu/BdSLT.pdf>
- [3] J.-D. Boissonnat, L. J. Guibas, and S. Y. Oudot. Manifold reconstruction in arbitrary dimensions using witness complexes. *Proc. 23rd Ann. Sympos. Comput. Geom.* (2007), 194–203.
- [4] M. A. Buchner. Simplicial structure of the real analytic cut locus. *Proc. American Math. Soc.* **64** (1977), 118–121.
- [5] F. Chazal and A. Lieutier. Topology guaranteeing manifold reconstruction using distance function to noisy data. *Proc. 22nd Ann. Sympos. Comput. Geom.* (2006), 112–118.
- [6] F. Chazal and A. Lieutier. The λ -medial axis. *Graphical Models* **67** (2006), 304–331.
- [7] F. Chazal and S. Oudot. Towards persistence-based reconstruction in Euclidean spaces. *Proc. 24th Ann. Sympos. Comput. Geom.* (2008), 232–241.

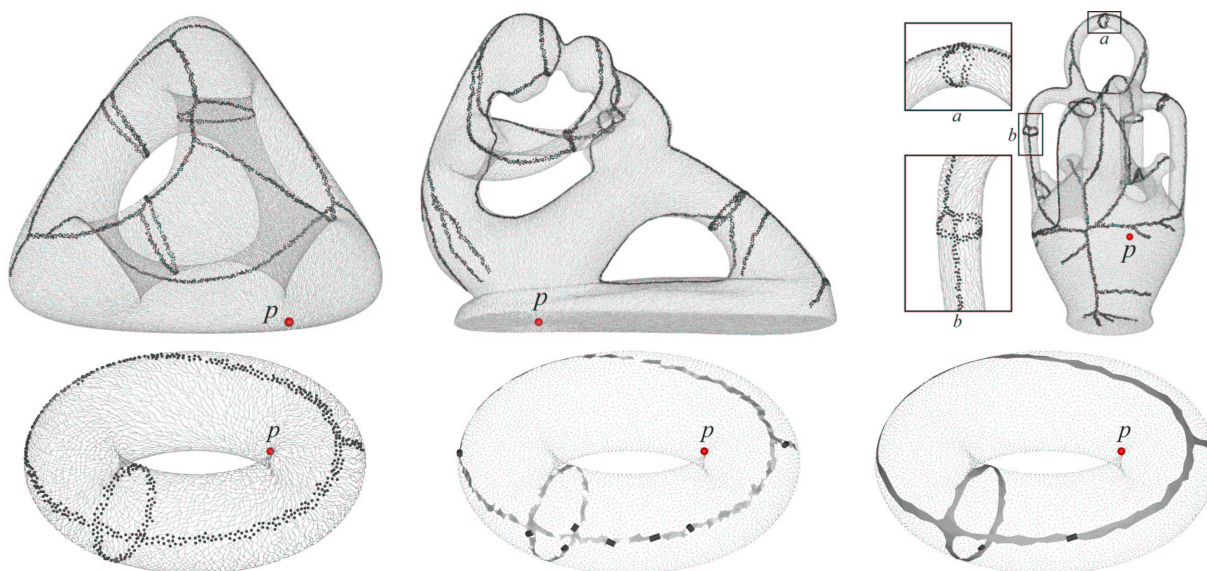


Figure 4: Top row: Cut locus approximation. Bottom row: Rips complexes of points approximating a cut locus in Torus; 2 out of 8 unpaired edges (shown dark) computed by the persistence algorithm in the middle Rips complex persist in the larger Rips complex on right since $\beta_1 = 2$ here.

- [8] E. W. Chambers, J. Erickson, and P. Worah. Testing contractibility in planar Rips complexes *Proc. 24th Ann. ACM Sympos. Comput. Geom.* (2008), 251–259.
- [9] I. Chavel. *Riemannian Geometry: A Modern Introduction*. Cambridge U. Press, New York, 1994.
- [10] S.-W. Cheng, T. K. Dey, and E. A. Ramos. Manifold reconstruction from point samples. *Proc. 16th Sympos. Discrete Algorithms* (2005), 1018–1027.
- [11] D. Cohen-Steiner, H. Edelsbrunner, and J. Harer. Stability of persistence diagrams. *Discrete Comput. Geom.* **37** (2007), 103–120.
- [12] É. Colin de Verdière and J. Erickson. Tightening non-simple paths and cycles on surfaces. *Proc. ACM-SIAM Sympos. Discrete Algorithms* (2006), 192–201.
- [13] T. K. Dey and K. Li. Topology from data via geodesic complexes. *Tech report OSU-CISRC-3/09-TR05*, 2009.
- [14] M. P. do Carmo. *Riemannian geometry*. Birkhäuser, Boston, 1992.
- [15] H. Edelsbrunner, D. Letscher, and A. Zomorodian. Topological persistence and simplification. *Discrete Comput. Geom.* **28** (2002), 511–533.
- [16] J. Erickson and K. Whittlesey. Greedy optimal homotopy and homology generators. *Proc. ACM-SIAM Sympos. Discrete Algorithms* (2005), 1038–1046.
- [17] H. Federer. Curvature measures. *Trans. Amer. Math. Soc.* **93** (1959), 418.
- [18] J. Gao, L. J. Guibas, S. Y. Oudot, and Y. Wang. Geodesic Delaunay triangulations and witness complexes in the plane. *Proc. ACM-SIAM Sympos. Discrete Algorithms* (2008), 571–580.
- [19] K. Hildebrandt, K. Polthier, and M. Wardetzky. On the convergence of metric and geometric properties of polyhedral surfaces. *Geometriae Dedicata* **123** (2006), 89–112.
- [20] I. T. Jolliffe. *Principal Component Analysis*. Springer series in statistics, Springer, NY, 2002.
- [21] W. P. A. Klingenberg. *Riemannian Geometry*. Walter de Gruyter, Berlin, 1995.

- [22] S. B. Myers. Connections between differential geometry and topology I: Simply connected surfaces. *Duke Math. J.* **1** (1935), 376–391.
- [23] S. B. Myers. Connections between differential geometry and topology II: Closed surfaces. *Duke Math. J.* **2** (1936), 95–102.
- [24] P. Niyogi, S. Smale, and S. Weinberger. Finding the homology of submanifolds with high confidence from random samples. *Discrete Comput. Geom.* (2006).
- [25] V. de Silva and G. Carlsson. Topological estimation using witness complexes. *Proc. Sympos. Point-Based Graph.* (2004), 157–166.
- [26] J. B. Tenenbaum, V. de Silva, and J. C. Langford. A global geometric framework for nonlinear dimensionality reduction. *Science* **290** (2000).
- [27] A. Zomorodian and G. Carlsson. Computing persistent homology. *Discr. Comput. Geom.* **33** (2005), 249–274.