

JPLEX WITH MATLAB TUTORIAL

HENRY ADAMS
JULY 19, 2009

CONTENTS

1. Introduction	1
1.1. JPlex	1
1.2. Accompanying files	2
1.3. Installation for MATLAB	2
1.4. Advanced options	3
2. Math review	3
2.1. Simplicial complexes	3
2.2. Homology	3
2.3. Filtered simplicial complexes	4
2.4. Persistent homology	4
3. Streams	4
3.1. Class SimplexStream	4
3.2. Subclass ExplicitStream and homology	4
3.3. Subclass ExplicitStream and persistent homology	6
3.4. ExplicitStream details	7
4. Point cloud data	8
4.1. Class PointData	8
4.2. Subclass EuclideanArrayData	8
4.3. Subclass DistanceData	9
5. Streams from point cloud data	10
5.1. Subclass RipsStream	10
5.2. Landmark selection	12
5.3. Subclass WitnessStream	14
5.4. Subclass LazyWitnessStream	15
6. Example with real data	16
Appendices	18
Appendix A. Dense core subsets	18
References	20

1. INTRODUCTION

1.1. **JPlex.** JPlex is a Java software package for computing the persistent homology of filtered simplicial complexes, often generated from point cloud data. The authors are Harlan Sexton and Mikael Vejdemo Johansson. This is a tutorial for using JPlex with MATLAB. There is a similar document for using JPlex with BeanShell. Please email Henry Adams (henrya@stanford.edu) with questions

about this tutorial or with suggestions for how to improve it. Some of the exercises are borrowed from Vin de Silva's *Plexercises*. Other sources of information about JPlex are the javadoc tree for the JPlex library (<http://comptop.stanford.edu/programs/jplex/files/javadoc/index.html>) and the TMSCHS website (<http://comptop.stanford.edu/programs/>).

1.2. Accompanying files. This tutorial should be accompanied by the following files.

- (1) `commandListMatlab.rtf`
- (2) `coreSubset.m`
- (3) `dct.m`
- (4) `exerciseAnswersMatlab.rtf`
- (5) `kDensityList.m`
- (6) `maxminLandmarks.m`
- (7) `PlexMatlabTutorial.pdf`
- (8) `pointsFigure8.m`
- (9) `pointsRange.mat`
- (10) `pointsTorus.m`
- (11) `startJPlex.m`

PDF file (7) is this tutorial. Text file (1) lists all commands in this tutorial, so that you may copy and paste long commands into the MATLAB window. Text file (4) has answers for the exercises in this tutorial.

The remaining files are MATLAB files. M-file (11) is a script: it is simply a shortcut that calls a sequence of MATLAB commands. M-files (2), (3), (5), (6), (8), and (10) are functions: they accept input and produce output. Writing your own such scripts or functions can be very useful. MAT-file (9) contains data.

1.3. Installation for MATLAB. Open MATLAB and check which version of Java is being used. In this tutorial, the symbol `>>` precedes commands to enter into your MATLAB window.

```
>> version -java
ans = Java 1.5.0.13 with Apple Inc. Java Hotspot(TM) Client
      VM mixed mode, sharing
```

JPlex requires version number 1.5 or higher.

Copy MATLAB files (2), (3), (5), (6), (8), (9), (10), and (11) from §1.2 into the current MATLAB directory.

Copy `plex.jar` into a directory of your choice. If you choose the current MATLAB directory, then you do not need to edit m-file `startJPlex.m` and can skip the rest of this paragraph. Suppose instead you have chosen the directory `/Users/myName/`. Open the m-file `startJPlex.m`, which contains the following lines.

```
javaaddpath('plex.jar')
import edu.stanford.math.plex.*;
```

Replace 'plex.jar' above with '/Users/myName/plex.jar' (substitute the path of your chosen directory). Save `startJPlex.m`.

Run the `startJPlex.m` file.

```
>> startJPlex
```

Installation is complete. Confirm that JPlex is working properly with the following command.

```
>> Simplex.makePoint(1,2)
ans = <(2) 1> % You should get this answer
```

Each time upon starting a new MATLAB session, you will need to run `startJPlex.m`.

1.4. Advanced options. Depending on the size of your JPlex computations, you may need to increase the maximum Java heap size. This will not be necessary for the examples in this tutorial. It has also not been necessary for my research with JPlex.

The following command returns your maximum heap size in bytes.

```
>> java.lang.Runtime.getRuntime.maxMemory
ans = 130875392
```

My computer has a heap limit of approximately 128 megabytes. To increase your limit to, say, 256 megabytes, create a file named `java.opts` in your MATLAB directory which contains the text `-Xmx256m` and then restart MATLAB.

2. MATH REVIEW

Below is a brief math review. For more details, see [2, 6, 9].

2.1. Simplicial complexes. An abstract simplicial complex is given by the following data.

- A set Z of vertices or 0-simplices.
- For each $k \geq 1$, a set of k -simplices $\sigma = [z_0 z_1 \dots z_k]$, where $z_i \in Z$.
- Each k -simplex has $k+1$ faces obtained by deleting one of the vertices. The following membership property must be satisfied: if σ is in the simplicial complex, then all faces of σ must be in the simplicial complex.

We think of 0-simplices as vertices, 1-simplices as edges, 2-simplices as triangular faces, and 3-simplices as tetrahedrons.

2.2. Homology. Betti numbers help describe the homology of a simplicial complex X . The value $Betti_k$, where $k \in \mathbb{N}$, is equal to the rank of the k -th homology group of X . Roughly speaking, $Betti_k$ gives the number of k -dimensional holes. In particular, $Betti_0$ is the number of connected components. For instance, a k -dimensional sphere has all Betti numbers equal to zero except for $Betti_0 = Betti_k = 1$.

2.3. Filtered simplicial complexes. A filtration on a simplicial complex X is a collection of subcomplexes $\{X(t) \mid t \in \mathbb{R}\}$ of X such that $X(t) \subset X(s)$ whenever $t \leq s$. The filtration time of a simplex $\sigma \in X$ is the smallest t such that $\sigma \in X(t)$. In JPlex, filtered simplicial complexes are called streams.

2.4. Persistent homology. Betti intervals help describe how the homology of $X(t)$ changes with t . A k -dimensional Betti interval, with endpoints $[t_{start}, t_{end})$, corresponds roughly to a k -dimensional hole that appears at filtration time t_{start} , remains open for $t_{start} \leq t < t_{end}$, and closes at time t_{end} . We are often interested in Betti intervals that persist for a long filtration range.

Persistent homology depends heavily on functoriality: for $t \leq s$, the inclusion $i : X(t) \rightarrow X(s)$ of simplicial complexes induces a map $i_* : H_k(X(t)) \rightarrow H_k(X(s))$ between homology groups.

3. STREAMS

3.1. Class SimplexStream. In JPlex, a filtered simplicial complex is called a stream, and streams are implemented by the class SimplexStream. The subclass ExplicitStream allows us to build a SimplexStream instance from scratch. In §5 we will learn about subclasses RipsStream, WitnessStream, and LazyWitnessStream, which construct SimplexStream instances from point cloud data.

3.2. Subclass ExplicitStream and homology. Since JPlex is designed to compute persistent homology, it is not the most efficient software (notationally or computationally) for computing homology.

Circle example. Let's build a simplicial complex homeomorphic to a circle. To build a simplicial complex in JPlex we simply build a stream in which all filtration times are zero. First we get an empty ExplicitStream instance.

```
>> s1=ExplicitStream;
```

Many command lines in this tutorial will end with a semicolon to suppress unwanted output such as

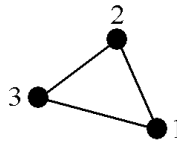
```
>> s1=ExplicitStream
s1 = edu.stanford.math.plex.ExplicitStream@b138fc
```

Next we add simplices using the method `add`. Two inputs are required: a matrix whose rows contain the simplices to be added, and a vector listing the corresponding filtration times. We choose all filtration times to be zero as we are building a simplicial complex instead of a stream.

```
>> s1.add([1;2;3], [0;0;0])           % adds 3 vertices
>> s1.add([1,2;2,3;3,1], [0;0;0])     % adds 3 edges
```

Let's inspect what we've built. The object `s1.dump(k)` contains information about the k -simplices of `s1`. We display the 1-simplices and their filtration times.

```
>> s1.dump(1).C                       % C stands for cell
ans =
```



```

      1  2
      1  3
      2  3
>> s1.dump(1).F           % F stands for filtration
ans =
      0
      0
      0

```

Our `s1`, like any `ExplicitStream` instance, has two states: open or closed. The state is automatically set to open whenever we edit or view the simplices of `s1`. We must manually close `s1` before computing its homology.

```
>> s1.close
```

The strange command names for computing homology will make sense later when we compute persistent homology.

```

>> intervals=Plex.Persistence.computeIntervals(s1);
>> Plex.FilterInfinite(intervals)
ans = BN{1, 1}

```

The result `BN{1, 1}` means that `s1` has $Betti_0 = 1$ and $Betti_1 = 1$, which are the Betti numbers of a circle. (If you instead get `BN{1}`, email Henry to get the most recently updated `plex.jar` file.)

6-sphere example. Let's build a 6-sphere, which is homeomorphic to the boundary of a 7-simplex. Adding the simplices manually as we did in the circle case would be very tedious: there are 8 vertices, 28 edges, 56 triangles, etc. Instead, we start with a 7-simplex that has no faces. We add its faces using the method `ensure_all_faces`. Then, we remove the 7-simplex, leaving only its boundary.

```

>> s6=ExplicitStream;
>> s6.add(1:8,0)           % 1:8 is the matrix [1,2,...,8]
                           % adds the 7-simplex

>> s6.ensure_all_faces
>> s6.remove(1:8)

```

How many 3-simplices did we avoid adding manually?

```

>> size(s6.dump(3).C)
ans = 70    4           % 70 3-simplices

```

We compute the homology.

```

>> s6.close
>> intervals=Plex.Persistence.computeIntervals(s6);
>> Plex.FilterInfinite(intervals)
ans = BN{1, 0, 0, 0, 0, 0, 1}

```

We get nonzero Betti numbers $Betti_0 = Betti_6 = 1$. (If you instead get `BN{1, 0, 0, 0, 0, 0}`, email Henry to get the most recently updated `plex.jar` file.)

The following command tells us that we have been computing homology over the coefficient field \mathbb{Z}_{11} .

```
>> Persistence.baseModulus
ans = 11
```

We can instead compute over modulus 13 (or any other prime between 2 and 251).

```
>> Persistence.setBaseModulus(13)
```

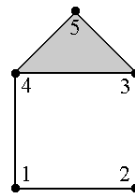
Exercise 3.2.1. Build a simplicial complex homeomorphic to the torus. Compute its Betti numbers. *Hint:* You will need at least 7 vertices [6, page 107]. I recommend using a 3×3 grid of 9 vertices.

Exercise 3.2.2. Build a simplicial complex homeomorphic to the Klein bottle. Check that it has the same Betti numbers as the torus over \mathbb{Z}_2 coefficients but different Betti numbers over \mathbb{Z}_3 coefficients.

Exercise 3.2.3. Build a simplicial complex homeomorphic to the projective plane. Compare its Betti numbers to those of the Klein bottle over \mathbb{Z}_2 and \mathbb{Z}_3 coefficients.

3.3. Subclass `ExplicitStream` and persistent homology.

Let's build a stream with nontrivial filtration times. We build a house, with the square appearing at time 0, the top vertex at time 1, the roof edges at times 2 and 3, and the roof 2-simplex at time 7.



```
>> house=ExplicitStream;
>> house.add([1;2;3;4;5], [0;0;0;0;1])
>> house.add([1,2;2,3;3,4;4,1;3,5;4,5], [0;0;0;0;2;3])
>> house.add([3,4,5], 7)
```

We compute the Betti intervals.

```
>> house.close
>> intervals=Plex.Persistence.computeIntervals(house);
```

There are four intervals.

```
>> length(intervals)
ans = 4
```

The fourth interval is a $Betti_1$ interval, starting at filtration time 3 and ending at 7.

```
>> intervals(4).dimension
ans = 1
>> intervals(4).start
ans = 3
>> intervals(4).end
ans = 7
```

Or, we can display the fourth interval all at once.

```
>> intervals(4)
ans = [1: (3.000000, 7.000000)]
```

This 1-dimensional hole is formed by the three edges of the roof. It forms when edge [4, 5] appears at filtration time 3 and closes when 2-simplex [3, 4, 5] appears at

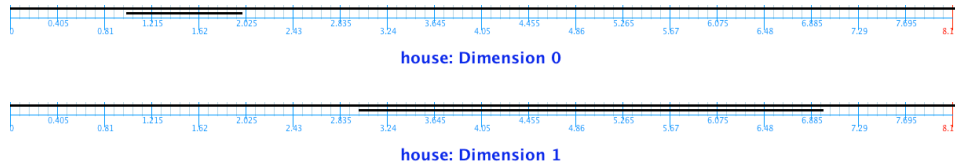
filtration time 7.

One $Betti_0$ interval and one $Betti_1$ interval are semi-infinite.

```
>> Plex.FilterInfinite(intervals)
ans = BN{1, 1}
```

The method `Plex.plot` lets us display the intervals as a Betti barcode. The three inputs are `intervals`, a string to appear as a label, and the maximum filtration time (which may be adjusted $\pm 10\%$) for the plot.

```
>> Plex.plot(intervals, 'house', 8)
```



The filtration times are on the horizontal axis. The $Betti_k$ number of the stream at filtration time t is the number of intervals in the dimension k plot that intersect a vertical line through t . Check that the displayed intervals agree with the filtration times we built into the stream `house`. At time 0, a connected component and a 1-dimensional hole form. At time 1, a second connected component appears, which joins to the first at time 2. A second 1-dimensional hole forms at time 3, and closes at time 7.

3.4. ExplicitStream details. We mention two remaining details about subclass `ExplicitStream`.

The methods `add` and `remove` do not necessarily enforce the definition of a stream. They allow us to build inconsistent streams in which some simplex $\sigma \in X(t)$ contains a subsimplex $\sigma' \notin X(t)$, meaning that $X(t)$ is not a simplicial complex. The method `verify(1)` returns 1 if our stream is consistent and returns 0 with explanation if not.

```
>> house.verify(1)
ans = 1
>> house.add([1,4,5],0)
>> house.verify(1)
Simplex <1, 4, 5> is present, but its face <1, 5> is not.
Simplex <1, 4, 5> has value 0.0000, but face <4, 5> has value
3.0000.
ans = 0
```

In §5 we will create `SimplexStream` instances that, unlike `house`, are not also `ExplicitStream` instances. To display or edit such streams, we first need to use the method `makeExplicit`. See Exercise 5.1.1.

4. POINT CLOUD DATA

4.1. Class `PointData`. A point cloud is a finite metric space. In JPLex, point cloud data is implemented by the class `PointData`. We detail two subclasses, `EuclideanArrayData` and `DistanceData`, that build `PointData` instances from different representations of a point cloud. In §5 we will learn how to build streams from a `PointData` instance.

4.2. Subclass `EuclideanArrayData`. This subclass is for a point cloud in a Euclidean space.

Let's give coordinates to the points of our house.

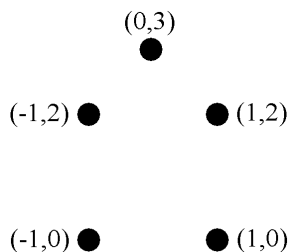


FIGURE 1. The house point cloud, stored in `PointData` instance `pdataHouse`

We create a `PointData` instance using these coordinates. The input to the `EuclideanArrayData` constructor is a matrix whose i -th row lists the coordinates of the i -th point.

```
>> pdataHouse=EuclideanArrayData([-1,0;1,0;1,2;-1,2;0,3]);
```

Any `PointData` instance can display the number of data points and the distance between, say, points 1 and 3.

```
>> pdataHouse.count
ans = 5
>> pdataHouse.distance(1,3)
ans = 2.8284
```

The method `dimension` returns the dimension of the Euclidean space containing our points.

```
>> pdataHouse.dimension
ans = 2
```

The method `coordinate(i,j)` returns the j -th coordinate of point i . Points are indexed starting at one but coordinates are indexed starting at zero. We display the coordinates of point 5.

```
>> pdataHouse.coordinate(5,0)
ans = 0
>> pdataHouse.coordinate(5,1)
ans = 3
```


The m-files `pointsFigure8.m` and `pointsTorus.m` each create a matrix of points which can be input into subclass `EuclideanArrayData`.

Figure 8 example. The first command below selects 100 points randomly from a figure eight (the union of unit circles centered at (0,1) and (0,-1)), and then adds noise to each point. The second command plots these points. The third command creates a `PointData` instance.

```
>> pointsF100=pointsFigure8(100);
>> plot(pointsF100(:,1),pointsF100(:,2),'.'), axis equal
>> pdataF100=EuclideanArrayData(pointsF100);
```

Torus example. The first command below selects 20^2 points from a 20×20 grid on the 2-dimensional unit torus in \mathbb{R}^4 , and then adds noise to each point. The second command plots the third and fourth coordinates of these points; the plot should be a circle. The third command creates a `PointData` instance.

```
>> pointsT20=pointsTorus(20);
>> plot(pointsT20(:,3),pointsT20(:,4),'.'), axis equal
>> pdataT20=EuclideanArrayData(pointsT20);
```

4.3. Subclass `DistanceData`. This subclass creates a `PointData` instance using a distance matrix. For a point cloud in Euclidean space, subclass `DistanceData` is generally less convenient than subclass `EuclideanArrayData`. However, subclass `DistanceData` can be used for a point cloud in an arbitrary metric space.

The matrix `distances` summarizes the metric for our house points in Figure 1: entry (i, j) is the distance from point i to point j . Don't forget you can copy and paste from `commandListMatlab.rtf` into the MATLAB window!

```
>> distances=[0,2,sqrt(8),2,sqrt(10);
2,0,2,sqrt(8),sqrt(10);
sqrt(8),2,0,2,sqrt(2);
2,sqrt(8),2,0,sqrt(2);
sqrt(10),sqrt(10),sqrt(2),sqrt(2),0]
```

```
distances =
```

0	2.0000	2.8284	2.0000	3.1623
2.0000	0	2.0000	2.8284	3.1623
2.8284	2.0000	0	2.0000	1.4142
2.0000	2.8482	2.0000	0	1.4142
3.1623	3.1623	1.4142	1.4142	0

We create a `PointData` instance from this matrix.

```
>> pdataHouseDD=DistanceData(distances);
```

Check that the methods `count` and `distance` return the same output with `pdataHouseDD` as with `pdataHouse`. The methods `dimension` and `coordinate` are

not functional with the DistanceData subclass.

5. STREAMS FROM POINT CLOUD DATA

In §3 we built instances of the class SimplexStream from scratch. In this section we construct streams from a point cloud Z . We use the three subclasses RipsStream, WitnessStream, and LazyWitnessStream, which build the Vietoris-Rips, witness, and lazy witness streams. See [4] for additional information.

All three subclasses take four of the same inputs: the granularity δ , the maximum dimension d_{max} , the maximum filtration time t_{max} , and a point cloud Z stored as a PointData instance. The first three inputs allow the user to limit the size of the constructed stream, for computational efficiency. No simplices above dimension d_{max} are included. The persistent homology of the resulting stream can be calculated only up to dimension $d_{max} - 1$ (do you see why?). Also, instead of computing complex $X(t)$ for all $t \geq 0$, we only compute $X(t)$ for $t = \delta, 2\delta, 3\delta, \dots, N\delta$, where N is the largest integer such that $N\delta \leq t_{max}$.

In this tutorial we use $\delta = 0.001$. If you ever choose d_{max} or t_{max} too large and MATLAB seems to be running forever, pressing the “control” and “c” buttons simultaneously sometimes halts the computation.

5.1. Subclass RipsStream. Let $d(\cdot, \cdot)$ denote the distance between two points. A natural stream to build is the Rips stream. The complex $\text{Rips}(Z, t)$ is defined as follows:

- the vertex set is Z .
- for vertices a and b , edge $[ab]$ is in $\text{Rips}(Z, t)$ if $d(a, b) \leq t$.
- a higher dimensional simplex is in $\text{Rips}(Z, t)$ if all of its edges are.

Note that $\text{Rips}(Z, t) \subset \text{Rips}(Z, s)$ whenever $t \leq s$, so the Rips stream is a filtered simplicial complex. Since a Rips complex is the maximal simplicial complex that can be built on top of its 1-skeleton, it is a *flag complex*.

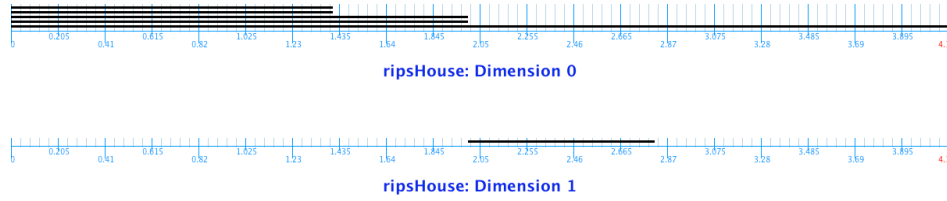
Let’s build a Rips stream instance `ripsHouse` from the PointData instance `pdataHouse`. Note this stream is different than the ExplicitStream `house` we built in §3.3.

```
>> ripsHouse=Plex.RipsStream(0.001,3,4,pdataHouse);
```

The order of the inputs is `RipsStream($\delta, d_{max}, t_{max}, Z$)`. Since $d_{max} = 3$ we can compute up to second dimensional persistent homology. For a Rips stream, the parameter t_{max} is the maximum possible edge length. Since $t_{max} = 4$ is greater than the diameter ($\sqrt{10}$) of our point cloud, all edges will eventually form.

We compute and display the Betti intervals. Typically the last input for the method `Plex.plot` will be t_{max} , since there is no reason to display filtration times that we haven’t computed.

```
>> intervals=Plex.Persistence.computeIntervals(ripsHouse);
>> Plex.plot(intervals,'ripsHouse',4)
```



The second dimensional Betti plot does not appear because there are no $Betti_2$ intervals. Check that these plots are consistent with the Rips definition: edges $[3, 5]$ and $[4, 5]$ appear at filtration time $t = \sqrt{2}$; the square appears at $t = 2$; the square closes at $t = \sqrt{8}$.

Exercise 5.1.1. Change `ripsHouse` into an explicit stream

```
>> ripsExpl=Plex.makeExplicit(ripsHouse);
```

Check that you can display and edit stream `ripsExpl` using the methods of §3.

Torus example. Try the following sequence of commands. We select a 20×20 grid of noisy points from a torus and build the RipsStream `ripsT20`. The fourth command returns the total number of simplices in `ripsT20`.

```
>> pointsT20=pointsTorus(20);
>> pdataT20=EuclideanArrayData(pointsT20);
>> ripsT20=Plex.RipsStream(0.001,3,0.9,pdataT20);
>> ripsT20.size
ans = 82831 % Generally close to 80000
>> intervals=Plex.Persistence.computeIntervals(ripsT20);
>> Plex.FilterInfinite(intervals)
ans = BN{1, 2, 1}
>> Plex.plot(intervals,'ripsT20',0.9)
```

The diameter of this torus (before adding noise) is $\sqrt{8}$, so choosing $t_{max} = 0.9$ likely will not show all homological activity. However, the torus will be reasonably connected by this time. Note the semi-infinite intervals match the correct numbers $Betti_0 = 1$, $Betti_1 = 2$, $Betti_2 = 1$ for a torus.

This example makes it clear that the computed “semi-infinite” intervals do not necessarily persist until $t = \infty$: in a Rips stream, once t is greater than the diameter of the point cloud, the Betti numbers for $\text{Rips}(Z, t)$ will be $Betti_0 = 1$, $Betti_1 = Betti_2 = \dots = 0$. The computed semi-infinite intervals are merely those that persist until $t = t_{max}$.

Exercise 5.1.2. Slowly increase the values for t_{max} , d_{max} , or the grid length (20 above) and note how quickly `ripsT20.size` and the computation time grows. Separately increasing t_{max} from 0.9 to 1, d_{max} from 3 to 4, or the grid length from 20 to 22 each roughly doubles `ripsT20.size`.

Exercise 5.1.3. Find a planar dataset Z and a filtration value t such that $Betti_2(\text{Rips}(Z, t)) \neq 0$. Build a RipsStream to confirm your answer.

Exercise 5.1.4. Find a planar dataset Z and a filtration value t such that $Betti_6(\text{Rips}(Z, t)) \neq 0$. When building a RipsStream to confirm your answer, don't forget to choose $d_{max} = 7$.

5.2. Landmark selection. For larger datasets, if we include every data point as a vertex, as in the Rips construction, our streams will quickly contain too many simplices for efficient computation. The witness stream and the lazy witness stream address this problem. In building these streams, we select a subset $L \subset Z$, called landmark points, as the only vertices. All data points in Z help serve as witnesses for the inclusion of higher dimensional simplices.

There are two common methods for selecting landmark points. The first is to choose the landmarks L randomly from point cloud Z . We select 25 random landmarks from figure eight PointData instance `pdataF100`.

```
>> L1=WitnessStream.makeRandomLandmarks(pdataF100,25);
>> length(L1)
ans = 26
```

Vector `L` always has first entry zero. The remaining 25 entries contain the indices of the random landmark vertices.

The second method for selecting landmark points, called sequential maxmin, is a greedy inductive selection process. Pick the first landmark randomly from Z . Inductively, if L_{i-1} is the set of the first $i-1$ landmarks, then let the i -th landmark be the point of Z which maximizes the function $z \mapsto d(z, L_{i-1})$, where $d(\cdot, \cdot)$ is the distance between the point and the set.

Landmarks chosen using sequential maxmin tend to cover the dataset and to be spread apart from each other. A disadvantage is that outlier points tend to be selected. Sequential maxmin landmarks are used in [1] and [3].

JPlex does not yet have a command for sequential maxmin landmark selection, so we use the m-file `maxminLandmarks.m`. The first input is a matrix and the second input is the number of landmarks to be selected. The third input is either the character `'e'` or `'d'`. Use `'e'` in the EuclideanArrayData case where the first input should be interpreted as an $N \times n$ matrix of N points in \mathbb{R}^n . Use `'d'` in the DistanceData case where the first input should be interpreted as an $N \times N$ distance matrix for N points in an arbitrary metric space.

Case EuclideanArrayData. Recall the figure eight example.

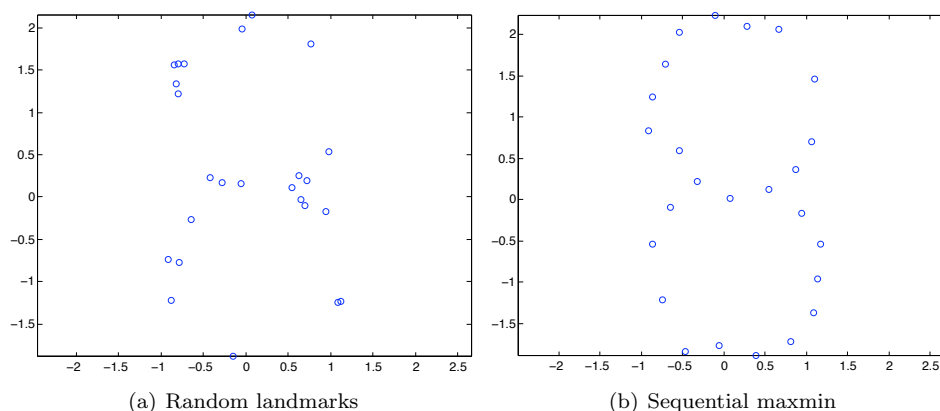
```
>> pdataF100=EuclideanArrayData(pointsF100);
```

We select 25 sequential maxmin landmarks.

```
>> L2=maxminLandmarks(pointsF100,25,'e');
```

Plot landmark points L1 and L2 to see the difference between random and sequential maxmin landmark selection.

```
>> pointsL1=pointsF100(L1(2:end),:);
>> plot(pointsL1(:,1),pointsL1(:,2),'o'), axis equal
>> pointsL2=pointsF100(L2(2:end),:);
>> figure; % This makes our next plot appear in a new window
>> plot(pointsL2(:,1),pointsL2(:,2),'o'), axis equal
```



Sequential maxmin seems to do a better job of choosing landmarks that cover the figure eight and that are spread apart.

Case DistanceData. Recall the house example.

```
>> pdataHouseDD=DistanceData(distances);
```

We select 3 sequential maxmin landmarks.

```
>> L=maxminLandmarks(distances,3,'d')
L =
```

```
0
3
1
2
```

Given point cloud Z and landmark subset L , we define $R = \max_{z \in Z} \{d(z, L)\}$. Number R reflects how finely the landmarks cover the dataset. We often use it as a guide for selecting the maximum filtration value t_{max} for a WitnessStream or LazyWitnessStream instance.

Exercise 5.2.1. Let Z be the point cloud in Figure 1 from §4.2, corresponding to PointData instance `pdataHouse`. Suppose we are using sequential maxmin to select a set L of 3 landmarks, and the first (randomly selected) landmark is $(1, 0)$. Find by hand the other two landmarks in L .

Exercise 5.2.2. Let Z be a point cloud and L a landmark subset. Show that if L is chosen via sequential maxmin, then for any $l_i, l_j \in L$, we have $d(l_i, l_j) \geq R$.

5.3. Subclass WitnessStream. Suppose we are given a point cloud Z and landmark subset L . Let $m_k(z)$ be the distance from a point $z \in Z$ to its $(k+1)$ -th closest landmark point. The witness stream complex $W(Z, L, t)$ is defined as follows.

- the vertex set is L .
- for $k > 0$ and vertices l_i , the k -simplex $[l_0 l_1 \dots l_k]$ is in $W(Z, L, t)$ if all of its faces are, and if there exists a witness point $z \in Z$ such that $\max\{d(l_0, z), d(l_1, z), \dots, d(l_k, z)\} \leq t + m_k(z)$.

Note that $W(Z, L, t) \subset W(Z, L, s)$ whenever $t \leq s$. Note that a landmark point can serve as a witness point.

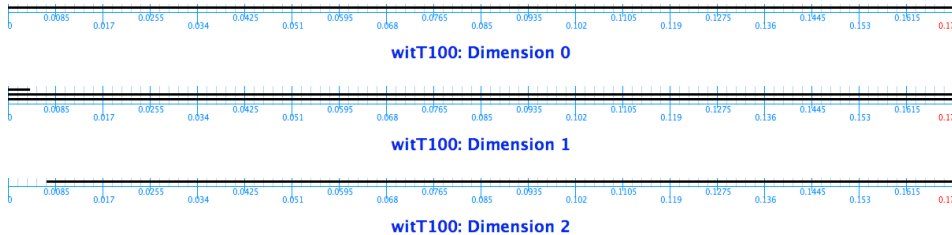
Exercise 5.3.1. Let Z be the point cloud in Figure 1 from §4.2, corresponding to PointData instance `pdataHouse`. Let $L = \{(1, 0), (0, 3), (-1, 0)\}$ be the landmark subset. Find by hand the filtration time for the edge between vertices $(1, 0)$ and $(0, 3)$. Which point or points witness this edge? What is the filtration time for the lone 2-simplex $[(1, 0), (0, 3), (-1, 0)]$?

Torus example. Let's build a WitnessStream instance for 100^2 points from a noisy torus, with 50 random landmarks. The fourth command returns the landmark covering measure R from §5.2. The fifth command returns our witness stream. The order of inputs is `WitnessStream($\delta, d_{max}, t_{max}, L, Z$)`. Often the value for t_{max} is chosen in proportion to R .

```
>> pointsT100=pointsTorus(100);
>> pdataT100=EuclideanArrayData(pointsT100);
>> L=WitnessStream.makeRandomLandmarks(pdataT100,50);
>> R=WitnessStream.estimateRmax(pdataT100,L)
R = 1.3188 % Generally close to 1.2
>> witT100=Plex.WitnessStream(0.001,3,R/8,L,pdataT100);
>> witT100.size
ans = 3320 % Generally close to 3000
```

We plot the Betti intervals.

```
>> intervals=Plex.Persistence.computeIntervals(witT100);
>> Plex.plot(intervals,'witT100',R/8)
```



The idea of persistent homology is that long intervals should correspond to real topological features, whereas short intervals are considered to be noise. The plot above shows that for a long range, the torus numbers $Betti_0 = 1$, $Betti_1 = 2$, $Betti_2 = 1$ are obtained. Your plot should contain a similar range.

The WitnessStream `witT100` contains approximately 3,000 simplices, fewer than the approximately 80,000 simplices in RipsStream `ripsT20`. This is despite the fact that we started with 100^2 points in the witness case, but only 20^2 points in the Rips case. This supports our belief that the witness stream returns good results at lower computational expense.

5.4. Subclass LazyWitnessStream. A lazy witness stream is similar to a witness stream. However, there is an extra parameter ν , typically chosen to be 0, 1, or 2, which helps determine how the lazy witness complexes $LW_\nu(Z, L, t)$ are constructed. See [4] for more information.

Suppose we are given a point cloud Z , landmark subset L , and parameter $\nu \in \mathbb{N}$. If $\nu = 0$, let $m(z) = 0$ for all $z \in Z$. If $\nu > 0$, let $m(z)$ be the distance from z to the ν -th closest landmark point. The lazy witness complex $LW_\nu(Z, L, t)$ is defined as follows.

- the vertex set is L .
- for vertices a and b , edge $[ab]$ is in $LW_\nu(Z, L, t)$ if there exists a witness $z \in Z$ such that $\max\{d(a, z), d(b, z)\} \leq t + m(z)$.
- a higher dimensional simplex is in $LW_\nu(Z, L, t)$ if all of its edges are.

Note that $LW_\nu(Z, L, t) \subset LW_\nu(Z, L, s)$ whenever $t \leq s$. The adjective *lazy* refers to the fact that the lazy witness complex is a flag complex: since the 1-skeleton determines all higher dimensional simplices, less computation is involved.

Exercise 5.4.1. Let Z be the point cloud in Figure 1 from §4.2, corresponding to PointData instance `pdataHouse`. Let $L = \{(1, 0), (0, 3), (-1, 0)\}$ be the landmark subset. Let $\nu = 1$. Find by hand the filtration time for the edge between vertices $(1, 0)$ and $(0, 3)$. Which point or points witness this edge? What is the filtration time for the lone 2-simplex $[(1, 0), (0, 3), (-1, 0)]$?

Exercise 5.4.2. Repeat the above exercise with $\nu = 0$ and with $\nu = 2$.

Exercise 5.4.3. Check that the 1-skeleton of a witness complex $W(Z, L, t)$ is the same as the 1-skeleton of a lazy witness complex $LW_2(Z, L, t)$. As a consequence, $LW_2(Z, L, t)$ is the flag complex of $W(Z, L, t)$.

The following sequence of commands is typical.

```
>> L=WitnessStream.makeRandomLandmarks(pdata,numLands);
           % Or, sequential maxmin landmarks
>> R=WitnessStream.estimateRmax(pdata,L);
>> laz=Plex.LazyWitnessStream(delta,d_max,t_max,nu,L,pdata);
```

```
>> intervals=Plex.Persistence.computeIntervals(laz);
>> Plex.plot(intervals,'laz',t_max)
```

Again, t_{max} is often chosen in proportion to R . In the next section we build a lazy witness stream on a dataset of range image patches.

6. EXAMPLE WITH REAL DATA

We now do an example with real data. Double check that the files `pointsRange.mat` and `dct.m`, which accompany this tutorial, are in your MATLAB directory.

In *On the nonlinear statistics of range image patches* [1], we study a space of range image patches drawn from the Brown database [7]. A range image is like an optical image, except that each pixel contains a distance instead of a grayscale value. Our space contains high-contrast, normalized, 5×5 pixel patches. We write each 5×5 patch as a length 25 vector and think of our patches as point cloud data in \mathbb{R}^{25} . We select from this space the 30% densest vectors, based on a density estimator called ρ_{300} (see Appendix A). In [1] this dense core subset is denoted $X^5(300,30)$, and it contains 15,000 points. In the next example we verify a result from [1]: $X^5(300,30)$ has the topology of a circle.

Load the file `pointsRange.mat`. The matrix `pointsRange` appears in your MATLAB workspace.

```
>> load pointsRange.mat
>> size(pointsRange)
ans = 15000    25           % 15000 points in dimension 25
```

Matrix `pointsRange` is in fact $X^5(300,30)$: each of its rows is a vector in \mathbb{R}^{25} . Display some of the coefficients of `pointsRange`. Can you visualize a circle?

We create a `PointData` instance using subclass `EuclideanArrayData`. We pick 50 sequential maxmin landmark points, find the value of R , and build the lazy witness stream with parameter $\nu = 1$.

```
>> pdataRange=EuclideanArrayData(pointsRange);
>> L=maxminLandmarks(pointsRange,50,'e');
>> R=WitnessStream.estimateRmax(pdataRange,L)
R = 0.7757           % Generally close to 0.75
>> lazRange=Plex.LazyWitnessStream(0.001,3,R/3,1,L,pdataRange);
>> lazRange.size
ans = 20937           % Generally between 10000 and 25000
>> intervals=Plex.Persistence.computeIntervals(lazRange);
>> Plex.plot(intervals,'lazRange',R/3)
```

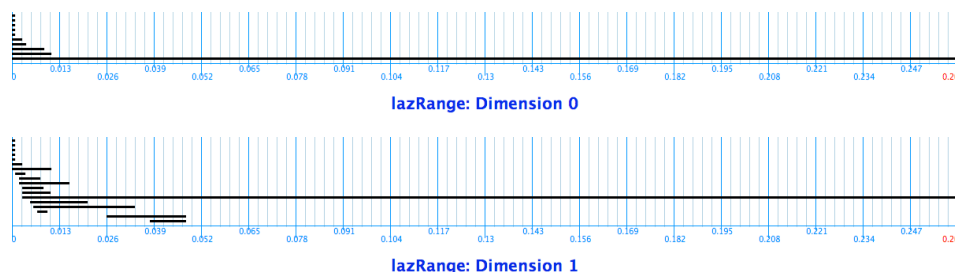



FIGURE 3. Betti intervals for `lazRange`, built from $X^5(300, 30)$

The plots above show that for a long range, the circle Betti numbers $Betti_0 = Betti_1 = 1$ are obtained. Your plot should contain a similar range. This is good evidence that the core subset $X^5(300, 30)$ is well-approximated by a circle.

Our 5×5 normalized patches are currently in the pixel basis: every coordinate corresponds to the range value at one of the 25 pixels. The Discrete Cosine Transform (DCT) basis is a useful basis for our patches [1, 7]. We change to this basis in order to plot a projection of the loop evidenced by Figure 3. The first command below computes the 5×5 DCT change-of-basis matrix. The second command changes basis via matrix multiplication.

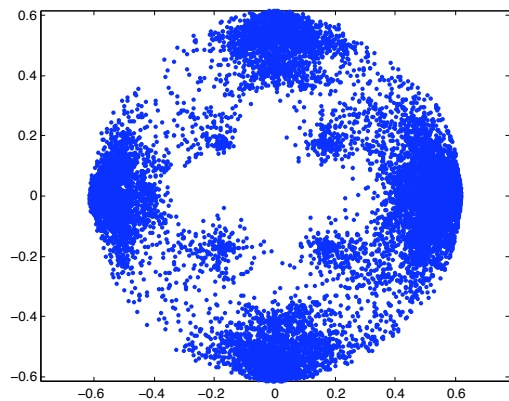
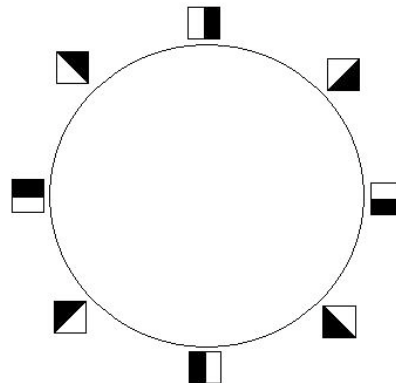
```
>> changeBasisDct=dct(5);
>> pointsRangeDct=pointsRange*changeBasisDct;
```

Two of the DCT basis vectors are horizontal and linear gradients.



We plot the projection of `pointsRangeDct` onto the linear gradient DCT basis vectors.

```
>> plot(pointsRangeDct(:,1),pointsRangeDct(:,5),'.'), axis
equal
```

(a) Projection of $X^5(300, 30)$ 

(b) Range primary circle

The projection of $X^5(300, 30)$ in Figure (a) shows a circle. It is called the range primary circle and is parameterized in Figure (b).

Appendices

APPENDIX A. DENSE CORE SUBSETS

A core subset of a dataset is a collection of the densest points, such as $X^5(300, 30)$ in §6. Since there are many density estimators, and since we can choose any number of the densest points, a dataset has a variety of core subsets. In this appendix we discuss how to create core subsets.

Real datasets can be very noisy, and outlier points can significantly alter the computed topology. Therefore, instead of trying to approximate the topology of an entire dataset, we often proceed as follows. We create a family of core subsets and identify their topologies. Looking at a variety of core subsets can give a good picture of the entire dataset.

See [3, 4] for an example using multiple core subsets. The dataset is high-contrast patches from natural images. The authors use three density estimators. As they change from the most global to the most local density estimate, the topologies of the core subsets change from a circle, to three intersecting circles, to a Klein bottle.

One way to estimate the density of a point z in a point cloud Z is as follows. Let $\rho_k(z)$ be the distance from z to its k -th closest neighbor. Let the density estimate at z be $\frac{1}{\rho_k(z)}$. Varying parameter k gives a family of density estimates. Using a small value for k gives a local density estimate, and using a larger value for k gives a more global estimate.

For Euclidean datasets, we use the m-file `kDensityList.m` to produce density estimates $\frac{1}{\rho_k}$. The following command is typical.

```
>> densities=kDensityList(points,k);
```

Input `points` is an $N \times n$ matrix of N points in \mathbb{R}^n . Input k is the density estimate parameter. Output `densities` is a vertical vector of length N containing the density estimate at each point.

M-file `coreSubset.m` builds a core subset. The following command is typical.

```
>> core=coreSubset(points,densities,numPoints);
```

Inputs `points` and `densities` are as above. Output `core` is a `numPoints` \times n matrix representing the `numPoints` densest points.

Prime numbers example. The command `primes(3571)` returns a vector listing all prime numbers less than or equal to 3571, which is the 500-th prime. We think of these primes as points in \mathbb{R} and build the core subset of the 10 densest points with density parameter $k = 1$.

```
>> p=primes(3571)';
>> length(p)
ans = 500
>> densities1=kDensityList(p,1);
>> core1=coreSubset(p,densities1,10)
core1 =
     2
     3
     5
     7
    11
    13
    17
    19
    29
    31
```

We get a bunch of twin primes, which makes sense since $k = 1$. Let's repeat with $k = 50$.

```
>> densities50=kDensityList(p,50);
>> core50=coreSubset(p,densities50,10)
core50 =
   113
   127
   109
   131
   107
   137
   139
   157
   149
   151
```

With $k = 50$, we expect the densest points to be slightly larger than the 25-th prime, which is 97.

Note: I typically use the m-file `kDensityList.m` on datasets of around 50,000 points (for instance, $X^5(300, 30)$ in §6 is the 30% densest points from a set of size 50,000, using the density estimate ρ_{300}). This computation takes about 10 minutes on my MacBook, which is much longer than any of the computations I do using the JPLex software.

REFERENCES

- [1] H. ADAMS AND G. CARLSSON, *On the nonlinear statistics of range image patches*, SIAM J. Img. Sci., 2, (2009), pp. 110–117.
- [2] M. A. ARMSTRONG, *Basic Topology*, Springer, New York, Berlin, 1983.
- [3] G. CARLSON, T. ISHKANOV, V. DE SILVA, AND A. ZOMORODIAN, *On the local behavior of spaces of natural images*, Int. J. Computer Vision, 76 (2008), pp. 1–12.
- [4] V. DE SILVA AND G. CARLSSON, *Topological estimation using witness complexes*, in Proceedings of the Symposium on Point-Based Graphics, ETH, Zürich, Switzerland, 2004, pp. 157–166.
- [5] H. EDELSBRUNNER, D. LETSCHER, AND A. ZOMORODIAN, *Topological persistence and simplification*, Discrete Computat. Geom., 28 (2002), pp. 511–533.
- [6] A. HATCHER, *Algebraic Topology*, Cambridge University Press, Cambridge, UK, 2002.
- [7] A. B. LEE, K. S. PEDERSEN, AND D. MUMFORD, *The nonlinear statistics of high-contrast patches in natural images*, Int. J. Computer Vision, 54 (2003), pp. 83–103.
- [8] H. SEXTON AND M. VEJDEMO-JOHANSSON, JPLex simplicial complex library. <http://comptop.stanford.edu/programs/jplex/>.
- [9] A. ZOMORODIAN AND G. CARLSSON, *Computing persistent homology*, Discrete Computat. Geom., 33 (2005), pp. 247–274.