# Topological Data Analysis and Persistence Theory

NSF/CBMS Conference

Valdosta State University

August 8-12, 2022

Peter Bubenik, University of Florida
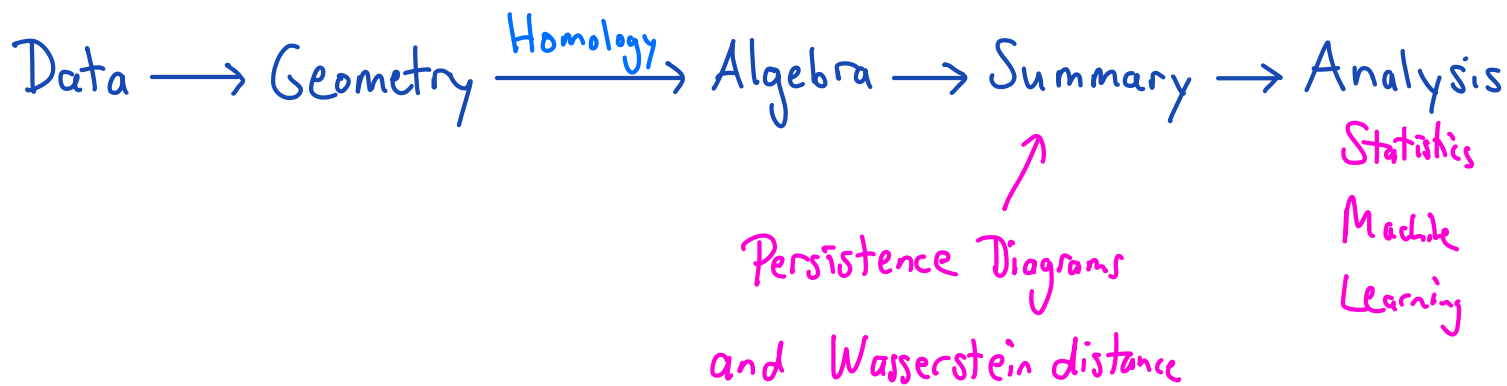
## Lecture 4: TDA and Statistics

Outline:
1. TDA and Hilbert Space
2. The Persistence Landscape
3. Statistics with Persistence Landscapes
4.

Please interrupt me!!!

# 1.  TDA and Hilbert Space

## 1.1  Why Hilbert Space?

Recall the TDA pipeline:

$$Data \longrightarrow Geometry \xrightarrow{Homology} Algebra \longrightarrow Summary \longrightarrow Analysis$$

Analysis: Statistics, Machine Learning

Summary ↑ Persistence Diagrams and Wasserstein distance

Statistics and Machine Learning depend on Linear Algebra.

Want: a vector space and inner product (ie inner product space).

We want summaries that lie in a complete inner product space (ie. a Hilbert space).

## 1.2  Nonembeddability

**Theorem**  For any $p \in [1, \infty]$, the metric space of persistence diagrams with the Wasserstein distance $W_p$ does not embed into a Hilbert space.

**Theorem**  For any $p \in (2, \infty]$, the metric space of persistence diagrams with the Wasserstein distance $W_p$ does not coarsely embed into a Hilbert space.

$p \in [1, 2]$   open

## 1.3  Feature maps and kernels

A **feature map** is a map $\Phi: X \longrightarrow H$.

<span style="color:magenta">Set         Hilbert space</span>

Given a feature map $\Phi$, we may define $k: X \times X \longrightarrow \mathbb{R}$

by $k(x, x') = \langle \Phi x, \Phi x' \rangle$, called a **kernel**.

<span style="color:magenta">inner product on $H$.</span>

**Theorem**   $k: X \times X \longrightarrow \mathbb{R}$ is a kernel iff it is symmetric and positive definite.

$\hookrightarrow \forall n \in \mathbb{N}, \forall x_1, \ldots x_n \in X,$ the $n \times n$ matrix $(k(x_i, x_j))$ has non-negative eigenvalues.
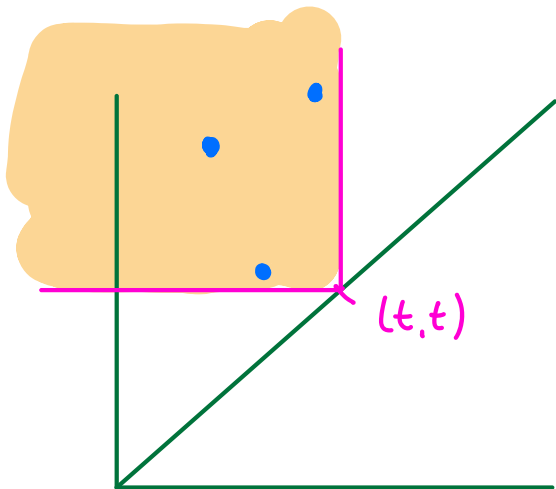
We want a Hilbert space $H$ and a feature map

$$\Phi: \mathrm{Dgm} \longrightarrow H.$$

# 2. The Persistence Landscape

## 2.1 Erosion of Persistence Diagrams

Let $M: \mathbb{R} \to \text{Vect}$ be a persistence module with persistence diagram $\text{Dgm } M = \{(b_i, d_i)\}_{i \in I}$
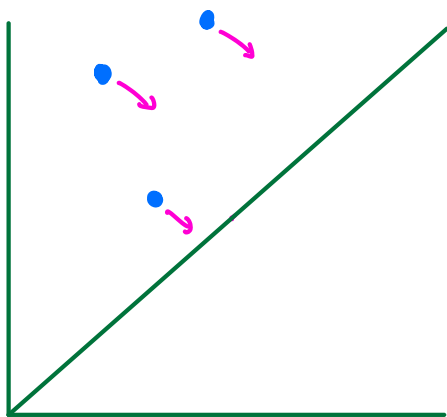


$(t, t)$

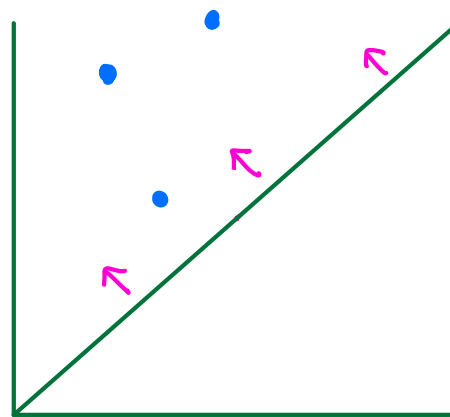**Lemma** $\dim M_t = \#$ of points in $\text{Dgm } M$ in the upper left quadrant $Q_t$.

Given $\varepsilon \geqslant 0$, let the $\underline{\varepsilon\text{-erosion}}$ of $\text{Dgm } M$ be given by $(\text{Dgm } M)_\varepsilon = \{(b_i + \varepsilon, d_i - \varepsilon)\}_{i \in I}$

$\{$ remove if $b_i + \varepsilon > d_i - \varepsilon$



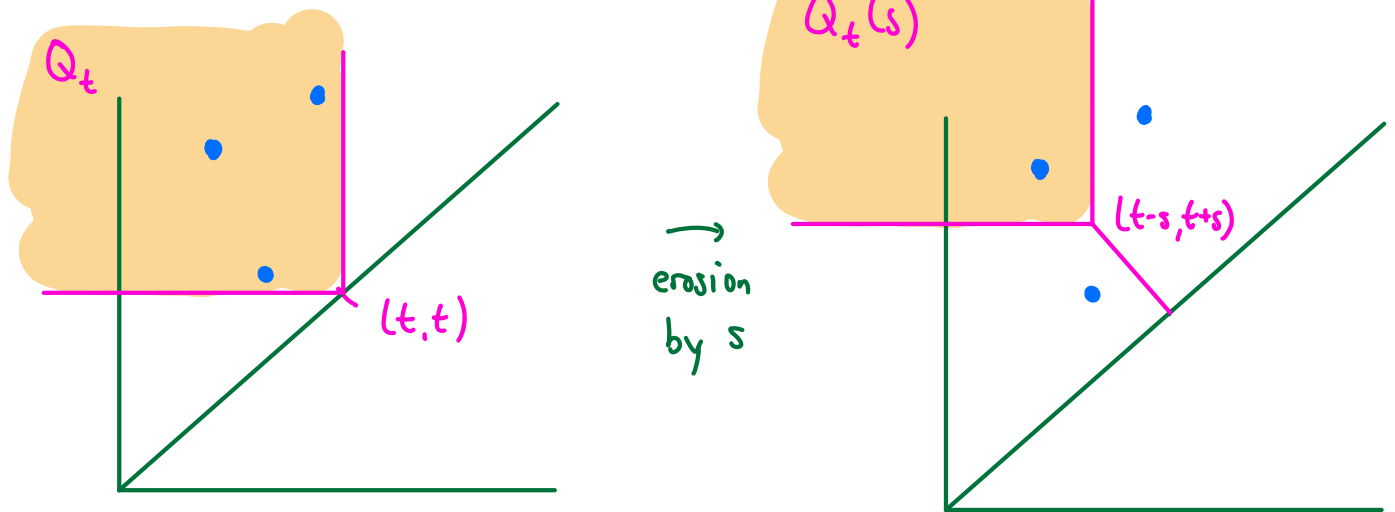OR



$\varepsilon$ confidence

# 2.2  The Persistence Landscape



Consider the # of points in the quadrant $Q_t(s)$.

For $k \in \mathbb{N}$,

Define $\lambda_k(t) = \max \left( s \mid \# \text{ points in } Q_t(s) \geq k \right)$

We obtain a sequence of functions $\lambda = (\lambda_1, \lambda_2, \lambda_3, \dots)$

called the <u>Persistence Landscape</u>.

It has inner product $\langle \lambda, \rho \rangle = \sum_k \int \lambda_k(t) \, \rho_k(t) \, dt$
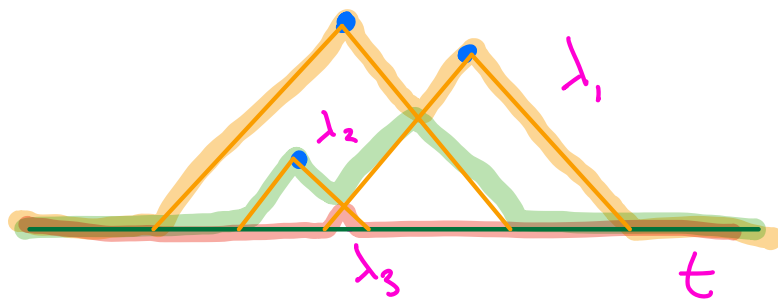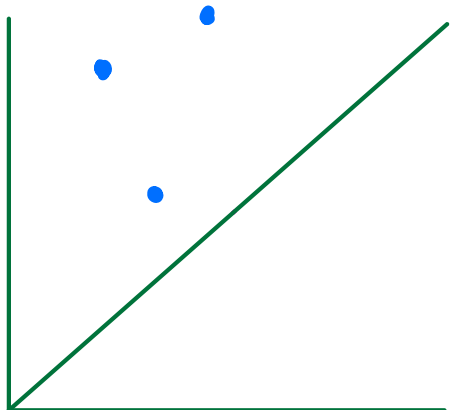
Equivalently, have $\lambda : \mathbb{N} \times \mathbb{R} \to \mathbb{R}$ given by $\lambda(k,t) = \lambda_k(t)$.

$\lambda \in L^2(\mathbb{N} \times \mathbb{R})$.

We have a feature map $\Lambda : \text{Dgm} \longrightarrow L^2(\mathbb{N} \times \mathbb{R})$

$$D \longmapsto \lambda$$

# Graphing the persistence landscape:
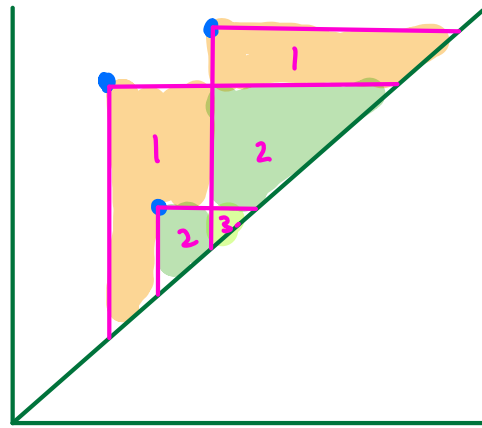


## Properties:

(a) PL. Each $\lambda_k$ is piecewise linear with slope $\pm 1$ on its support

(b) Lossless. Pers Diag $\rightarrow$ Pers Landscape is invertible

(c) Stable. Point cloud $\rightarrow$ Pers Landscape is nonexpansive

                           ↗                     ↑

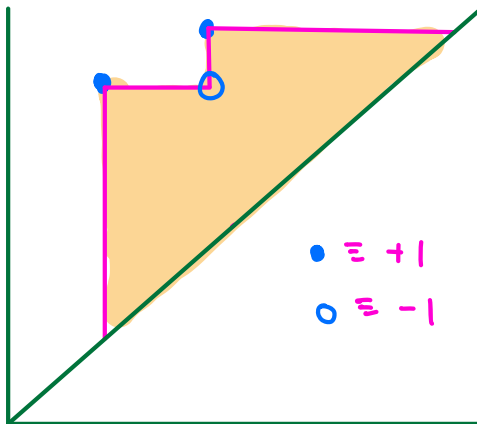                   Hausdoff dist         Supremum norm

## 2.3    Graded Persistence Diagram and Persistence Landscape

Rank function , Rank
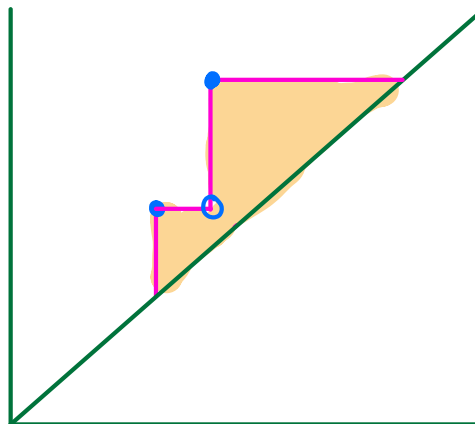and
Persistence Diagram, PD
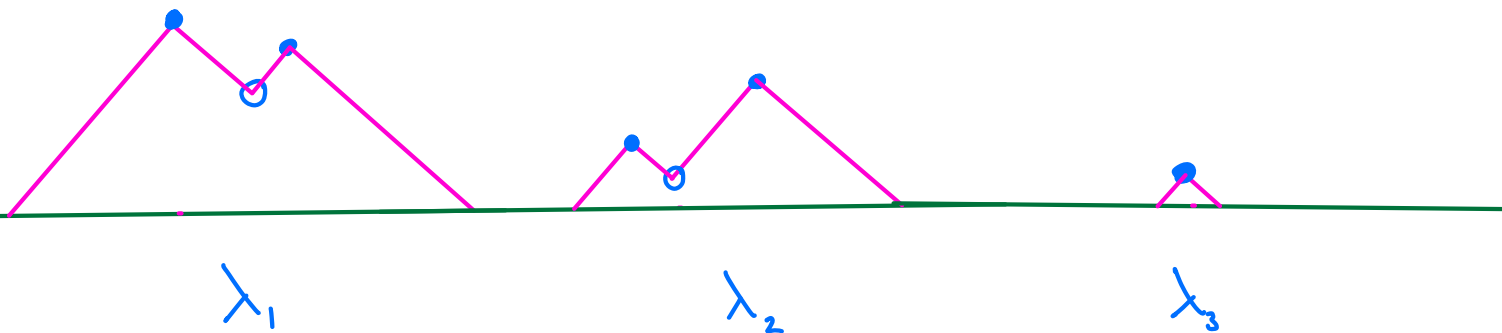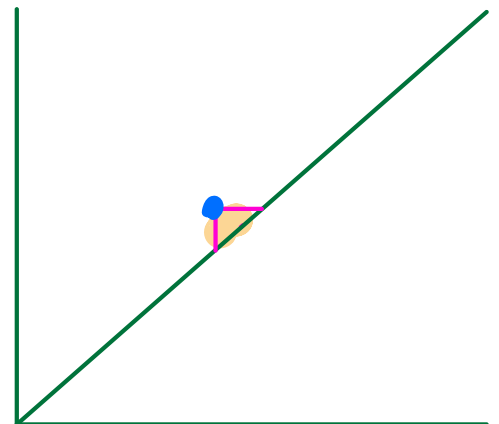
Graded Rank functions and Graded Persistence Diagrams:

Rank$_1$ and PD$_1$          Rank$_2$ and PD$_2$          Rank$_3$ and PD$_3$

$\bullet$ = +1
$\circ$ = -1

$\lambda_1$                    $\lambda_2$                    $\lambda_3$

## Theorem

The positive points of PD$_k$ are the local maxima of $\lambda_k$.
The negative points of PD$_k$ are the local minima of $\lambda_k$.

# 3. Statistics with Persistence Landscapes

## 3.1 Average Persistence Landscape

Data $\longrightarrow$ Persistence Landscape

$X$ $\qquad$ $\lambda$

Data $\longrightarrow$ Persistence Landscapes

$X_1, ..., X_N$ $\qquad$ $\lambda^{(1)}, ... \lambda^{(N)}$

Let $\bar{\lambda}(k,t) = \frac{1}{N} \lambda(k,t)$

## 3.2 Hypothesis Testing

If we have two experimental conditions

| Data | Average Persistence Landscapes |
|------|-------------------------------|

$X_1, ... X_N \longmapsto \bar{\lambda}$

$Y_1, ..., Y_N \longmapsto \bar{\rho}$

Difference of Average PL: $\bar{\lambda} - \bar{\rho}$

Is this difference significant?

Test statistic: $\| \bar{\lambda} - \bar{\rho} \|$

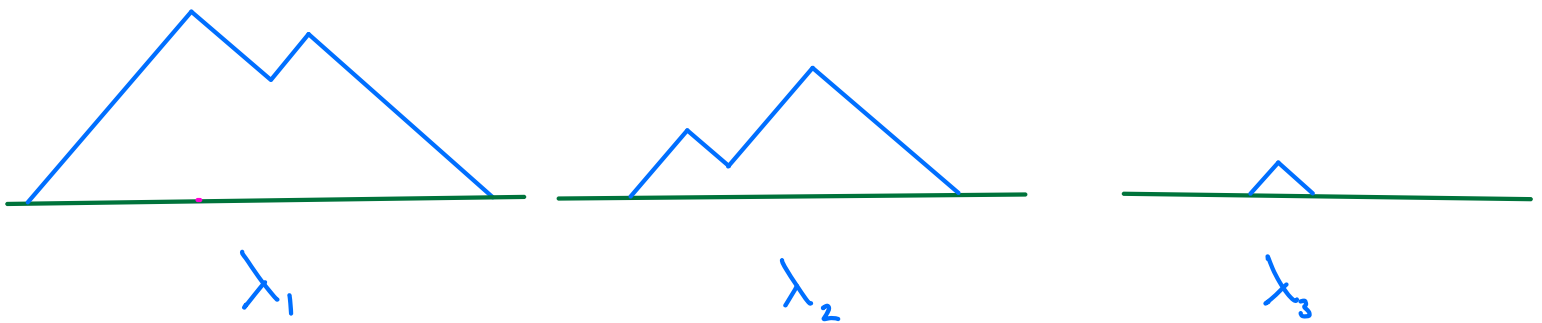We can obtain a p value for this statistic using a permutation test.

## 3.3 Discretizing the Persistence Landscape

$\lambda \in L^2(\mathbb{N} \times \mathbb{R})$ may be approximated by a point in $\mathbb{R}^D$ for some large $D$ as follows:

1. Discretize the support.
2. Evaluate $\lambda_1, \lambda_2, \ldots, \lambda_k$ on this grid.
3. Concatenate the numbers.



$\lambda_1 \qquad\qquad\qquad \lambda_2 \qquad\qquad\qquad \lambda_3$

## 3.4 Using Average Persistence Landscapes

If one is a situation where data is cheap and abundant then it is advised to repeat and use Average PL instead of PL.

If one is in a situation where the data is too large to compute persistent homology then subsample many times and compute the average persistence landscape.

**Theorem** In certain generic situations, not only can we use a persistence landscape to reconstruct a persistence diagram, we can use an average persistence landscape to reconstruct all of the persistence diagrams used to compute it.