A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

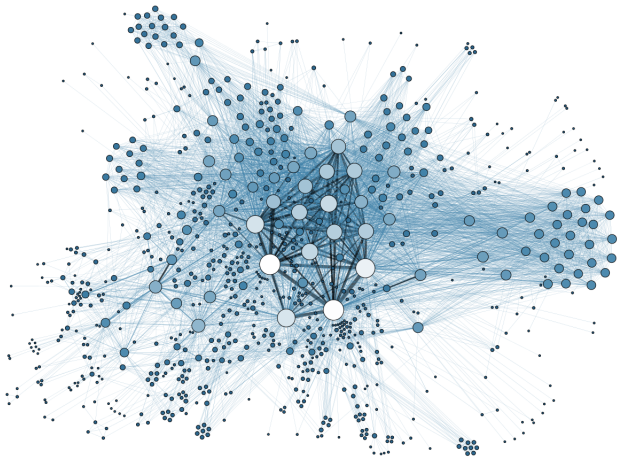# A Short Introduction to TDA (Topological Data Analysis)

Sara Kališnik

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Introduction

Data is often given in the form of point clouds in $\mathbb{R}^n$.

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Introduction

We have problems analyzing this data because it is often

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Introduction

We have problems analyzing this data because it is often

- given in the form of very long vectors, where not all coordinates are relevant,

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Introduction

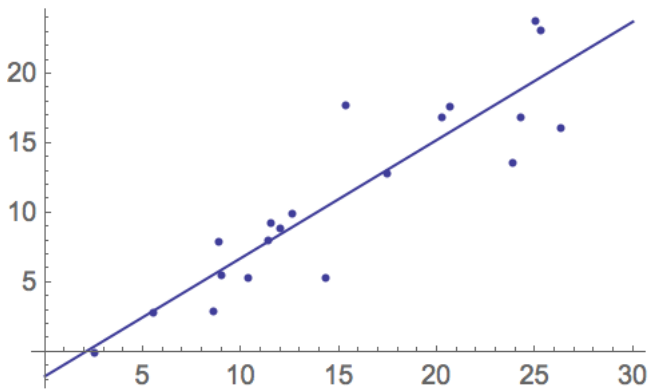We have problems analyzing this data because it is often

- given in the form of very long vectors, where not all coordinates are relevant,
- very high-dimensional,

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Introduction

We have problems analyzing this data because it is often

- given in the form of very long vectors, where not all coordinates are relevant,
- very high-dimensional,
- noisy.

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Introduction

We have problems analyzing this data because it is often

- given in the form of very long vectors, where not all coordinates are relevant,
- very high-dimensional,
- noisy.

Goal of topological data analysis:

Leverage machinery of algebraic topology to develop tools for studying 'qualitative' features of data.
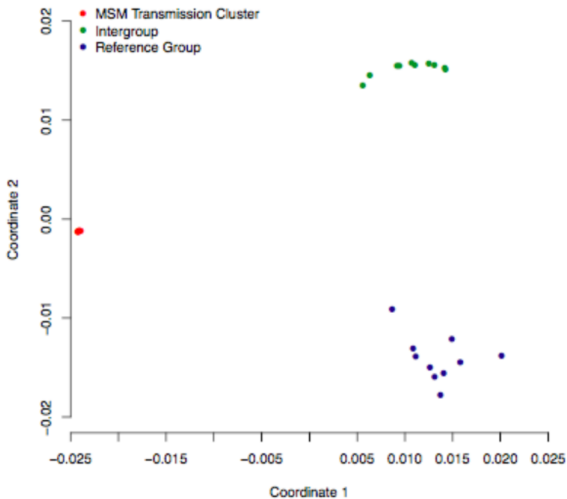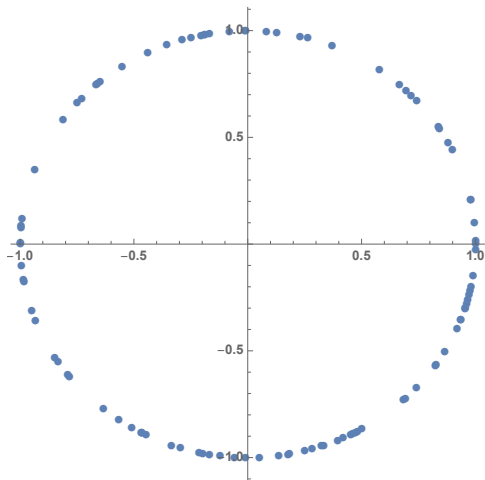
A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Shape of Data

## Linear Regression

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Shape of Data

## Clusters

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Shape of Data

## Loops

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Shape of Data

## Holes/Cycles/Loops

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Shape of Data

## Holes/Cycles/Loops

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Shape of Data

## Tendrils/Flares

Breast Cancer Study [Nicolau, Levine, Carlsson 2011]

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Why Topology?

Three key ideas:

- Invariance under deformation
- Coordinate freeness
- Compressed representations

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Why Topology?

Three key ideas:

- Invariance under deformation

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Why Topology?

Three key ideas:

- Invariance under deformation

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Why Topology?

Three key ideas:

- Coordinate Freeness

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Why Topology?

Three key ideas:

- Compressed representations

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# How to deal with shape?

Two tasks:

- Measure Shape
- Represent Shape

A Short
Introduction
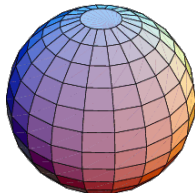to TDA
(Topological
Data Analysis)

Sara Kališnik

# Measuring Shape

Homology is a formalism for measuring shape...



$b_1 = 1$     $b_1 = 0$     $b_1 = 2$
$b_2 = 0$     $b_2 = 1$     $b_2 = 1$

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

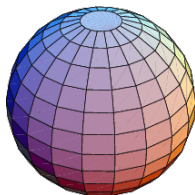# Measuring Shape

Homology is a formalism for measuring shape...
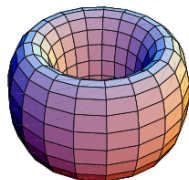


$b_1 = 1$       $b_1 = 0$       $b_1 = 2$
$b_2 = 0$       $b_2 = 1$       $b_2 = 1$

The extension of homology to more general setting including point clouds is called persistent homology.

The concept emerged independently in the work of Frosini, Ferri, and collaborators in Bologna, Italy, of Robins at Boulder, Colorado, and of Edelsbrunner, Letscher and Zomorodian at Duke, North Carolina.

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Persistent Homology

A finite metric space $\mathbb{X}$ has no interesting topology.

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Persistent Homology

A finite metric space $\mathbb{X}$ has no interesting topology.

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Persistent Homology

A finite metric space $\mathbb{X}$ has no interesting topology.



Let $U(\mathbb{X}, R)$ be the union of balls of radius $R$ centered at the points of $\mathbb{X}$. For any $R > 0$ and $i \geq 0$, $i$-th Betti number of $U(\mathbb{X}, R)$ gives us a qualitative descriptor of $\mathbb{X}$.

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Persistent Homology

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Persistent Homology



$b_0 = 1$
$b_1 = 2$

$b_0 = 1$
$b_1 = 1$

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Persistent Homology

Problems with this descriptor

- No canonical choice of $R$.

- Invariant is unstable with respect to perturbation of data or small changes in $R$.

- Does not distinguish 'small' holes from 'big' ones.

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Persistent Homology

## Persistent Homology

- Consider not only single reconstruction $U(\mathbb{X}, R)$ of $\mathbb{X}$, but a 1-parameter family of reconstructions

$$F(\mathbb{X}) = \{U(\mathbb{X}, r)\}_{r \in [0, \infty)}$$

and inclusion maps $U(\mathbb{X}, r) \hookrightarrow U(\mathbb{X}, r')$ whenever $r \leq r'$.

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Persistent Homology

## Persistent Homology

- Consider not only single reconstruction $U(\mathbb{X}, R)$ of $\mathbb{X}$, but a 1-parameter family of reconstructions

$$F(\mathbb{X}) = \{U(\mathbb{X}, r)\}_{r \in [0, \infty)}$$

  and inclusion maps $U(\mathbb{X}, r) \hookrightarrow U(\mathbb{X}, r')$ whenever $r \leq r'$.

- Apply $i$-dimensional homology functor $H_i$ with field coefficients

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Persistent Homology

## Persistent Homology

- Consider not only single reconstruction $U(\mathbb{X}, R)$ of $\mathbb{X}$, but a 1-parameter family of reconstructions

$$F(\mathbb{X}) = \{U(\mathbb{X}, r)\}_{r \in [0, \infty)}$$

and inclusion maps $U(\mathbb{X}, r) \hookrightarrow U(\mathbb{X}, r')$ whenever $r \leq r'$.

- Apply $i$-dimensional homology functor $H_i$ with field coefficients

- Obtain a family of vector spaces $\{V_r\}_r$ and linear maps between them. Call such algebraic structures persistence vector spaces.

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Persistent Homology

Persistent Homology

- Consider not only single reconstruction $U(\mathbb{X}, R)$ of $\mathbb{X}$, but a 1-parameter family of reconstructions

$$F(\mathbb{X}) = \{U(\mathbb{X}, r)\}_{r \in [0, \infty)}$$

  and inclusion maps $U(\mathbb{X}, r) \hookrightarrow U(\mathbb{X}, r')$ whenever $r \leq r'$.

- Apply $i$-dimensional homology functor $H_i$ with field coefficients

- Obtain a family of vector spaces $\{V_r\}_r$ and linear maps between them. Call such algebraic structures persistence vector spaces.

Can we classify persistence vector spaces that arise from filtrations up to isomorphism?

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Persistent Homology

Persistent Homology

- Consider not only single reconstruction $U(\mathbb{X}, R)$ of $\mathbb{X}$, but a 1-parameter family of reconstructions

$$F(\mathbb{X}) = \{U(\mathbb{X}, r)\}_{r \in [0, \infty)}$$

  and inclusion maps $U(\mathbb{X}, r) \hookrightarrow U(\mathbb{X}, r')$ whenever $r \leq r'$.
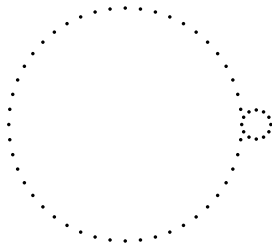
- Apply $i$-dimensional homology functor $H_i$ with field coefficients

- Obtain a family of vector spaces $\{V_r\}_r$ and linear maps between them. Call such algebraic structures persistence vector spaces.

Can we classify persistence vector spaces that arise from filtrations up to isomorphism?

Yes, by barcodes.
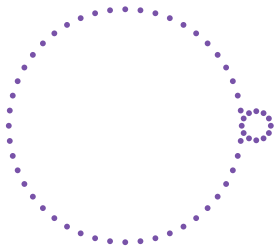
(Computing Persistent Homology, Carlsson and Zomorodian)

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Persistent Homology

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Persistent Homology

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Persistent Homology

Persistent Homology

A Short Introduction to TDA (Topological Data Analysis)

Sara Kališnik

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Persistent Homology

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Persistent Homology

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Persistent Homology

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Persistent Homology

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Persistent Homology

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Persistent Homology

Barcode for



$H_1$:

_____

_____

_____→

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Persistent Homology

Barcode for



$H_1$:

For each interval:

- Left endpoint is the index at which the hole is born
- Right endpoint is index at which hole dies
- Length of interval is the lifetime of a hole in filtration

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Applications of Persistent Homology

## Natural Scene Statistics/Image Processing

(Local structure of spaces of natural images by G. Carlsson, Vin de Silva, T. Ishkanov and A. Zomorodian)

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Natural Scene Statistics/Image Processing

A long time ago in a country far far away (the Netherlands) J. van Hateren and A. van der Schaaf were taking photos in a town called Groningen and in the surrounding countryside.

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Natural Scene Statistics/Image Processing

An image taken by black and white digital camera can be viewed as a vector, with one coordinate for each pixel.

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Natural Scene Statistics/Image Processing

An image taken by black and white digital camera can be viewed as a vector, with one coordinate for each pixel.

Typical camera uses tens of thousands of pixels, so images lie in a very high dimensional pixel space, $\mathbb{R}^P$.

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Natural Scene Statistics/Image Processing

An image taken by black and white digital camera can be viewed as a vector, with one coordinate for each pixel.

Typical camera uses tens of thousands of pixels, so images lie in a very high dimensional pixel space, $\mathbb{R}^P$.

David Mumford: What can be said about the set of images $\mathcal{I} \subseteq \mathcal{P}$ lying within $\mathbb{R}^P$? Can it be modeled as a submanifold or a subspace of $\mathbb{R}^P$?

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Natural Scene Statistics/Image
# Processing

The whole manifold of images is not accessible in a useful way,
a space of small image patches might in fact contain quite
useful information.

A Short
Introduction
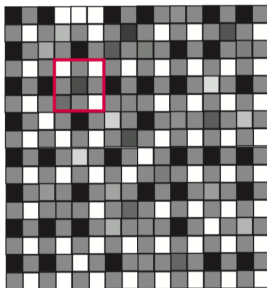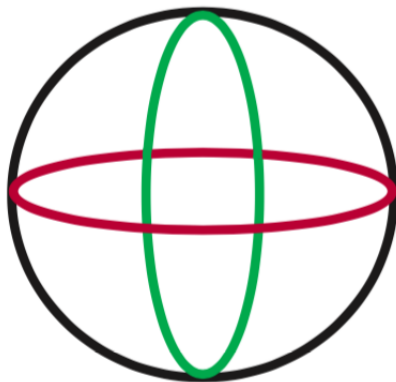to TDA
(Topological
Data Analysis)

Sara Kališnik

# Natural Scene Statistics/Image Processing

The whole manifold of images is not accessible in a useful way, a space of small image patches might in fact contain quite useful information.

Solution: observe $3 \times 3$ patches.

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Natural Scene Statistics/Image Processing



Three circle model

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Natural Scene Statistics/Image Processing



Three circle model in the data

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Natural Scene Statistics/Image Processing

Klein Bottle

**A Picture is worth 1,000 words**

J. Perea, G. Carlsson: Compression based on the Klein bottle mode (Kleinlets).

The evidence for Kleinlets over Wedglets

Original

Coded by Kleinlet at .71bpp
PSNR= 29dB

Coded by Wedgelet at .8bpp
PSNR= 27.7dB

Kleinlet     Wedgelet

Kleinlet

Wedgelet

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Applications of Persistent Homology

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Applications of Persistent Homology



Tree of Life

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Applications of Persistent Homology

- 1970s molecular phylogenetic analysis based on nucleotide and protein sequences

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Applications of Persistent Homology

- 1970s molecular phylogenetic analysis based on nucleotide and protein sequences
- 1977 Carl Woese identifies archaea as new domain in life

A Short
Introduction
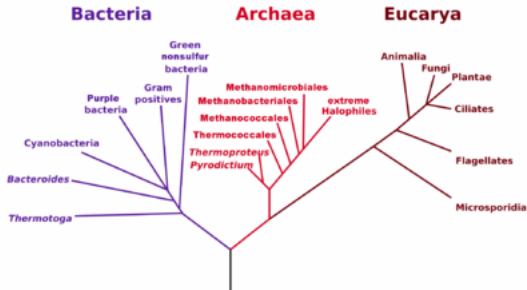to TDA
(Topological
Data Analysis)
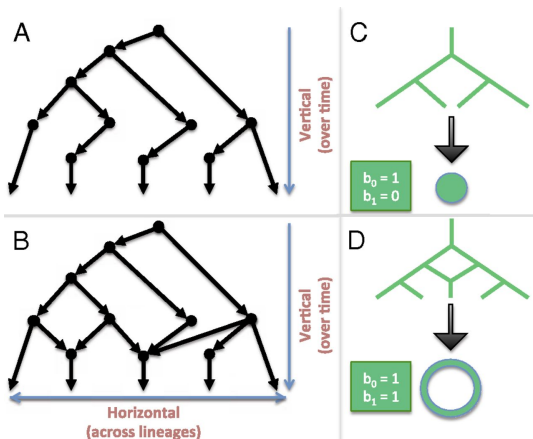
Sara Kališnik

# Applications of Persistent Homology

- 1970s molecular phylogenetic analysis based on nucleotide and protein sequences
- 1977 Carl Woese identifies archaea as new domain in life
- since 1990s a true revolution in genomic sequencing techniques providing hard data for evolutionary biology

**Phylogenetic Tree of Life**

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Applications of Persistent Homology

Viral Evolution (Topology of viral evolution by J.M. Chan, G. Carlsson, and R. Rabadan)

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Representing Shape

A very popular TDA method for representing shape is called mapper and was developed by G. Singh, F. Memoli and G. Carlsson.

A Short
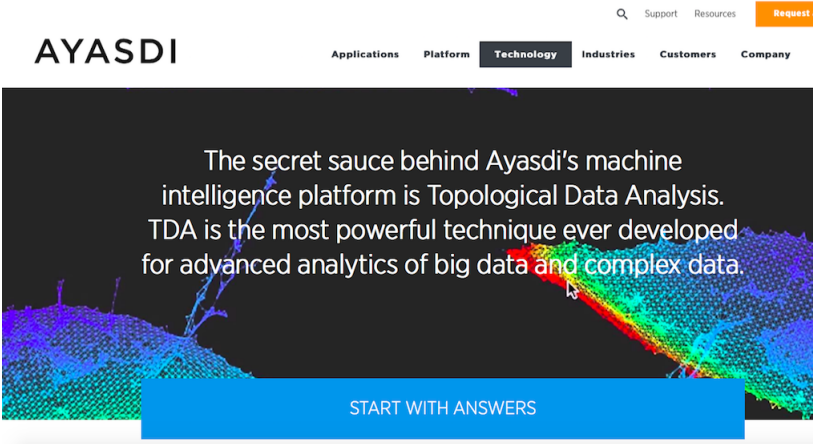Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Representing Shape

A very popular TDA method for representing shape is called mapper and was developed by G. Singh, F. Memoli and G. Carlsson.

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Representing Shape

Suppose we have a covering of a circle:

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Representing Shape

We assign a vertex to each connected component of this covering

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Representing Shape

When precisely two connected components intersect, we connect the corresponding vertices with an edge.

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Representing Shape

When precisely two connected components intersect, we connect the corresponding vertices with an edge.



When more than two, add a face of appropriate dimension.

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Representing Shape

Voila!

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Representing Shape

### Topological version of Mapper

Setting:

We are given a space $X$ equipped with a continuous map $f \colon X \to Z$ to a parameter space $Z$, and that the space $Z$ is equipped with a covering $\mathcal{U} = \{U_\alpha\}_{\alpha \in A}$ for some finite indexing set $A$.

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Representing Shape

### Topological version of Mapper

Setting:

We are given a space $X$ equipped with a continuous map
$f \colon X \to Z$ to a parameter space $Z$, and that the space $Z$ is
equipped with a covering $\mathcal{U} = \{U_\alpha\}_{\alpha \in A}$ for some finite
indexing set $A$.

- Since $f$ is continuous, the sets $f^{-1}(U_\alpha)$ form an open
  covering of $X$.

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Representing Shape

### Topological version of Mapper

Setting:

We are given a space $X$ equipped with a continuous map
$f: X \to Z$ to a parameter space $Z$, and that the space $Z$ is
equipped with a covering $\mathcal{U} = \{U_\alpha\}_{\alpha \in A}$ for some finite
indexing set $A$.

- Since $f$ is continuous, the sets $f^{-1}(U_\alpha)$ form an open
  covering of $X$.
- We write $f^{-1}(U_\alpha) = \cup_{j=1}^{j_\alpha} V(\alpha, i)$ where $j_\alpha$ is the number
  of connected components of $f^{-1}(U_\alpha)$. We write $\overline{\mathcal{U}}$ for the
  covering of $X$ obtained by taking these connected
  components.

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Representing Shape

### Topological version of Mapper

Setting:

We are given a space $X$ equipped with a continuous map
$f \colon X \to Z$ to a parameter space $Z$, and that the space $Z$ is
equipped with a covering $\mathcal{U} = \{U_\alpha\}_{\alpha \in A}$ for some finite
indexing set $A$.

- Since $f$ is continuous, the sets $f^{-1}(U_\alpha)$ form an open
  covering of $X$.
- We write $f^{-1}(U_\alpha) = \cup_{j=1}^{j_\alpha} V(\alpha, i)$ where $j_\alpha$ is the number
  of connected components of $f^{-1}(U_\alpha)$. We write $\overline{\mathcal{U}}$ for the
  covering of $X$ obtained by taking these connected
  components.
- Represent the topological space by a nerve of $\overline{\mathcal{U}}$.

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Representing Shape

The Statistical version of Mapper

- Define a reference map $f : X \to Z$, where $X$ is the given a point cloud and $Z$ is the reference metric space.

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Representing Shape

## The Statistical version of Mapper

- Define a reference map $f : X \to Z$, where $X$ is the given a point cloud and $Z$ is the reference metric space.

- Select a covering $\mathcal{U}$ of $Z$.

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Representing Shape

The Statistical version of Mapper

- Define a reference map $f : X \to Z$, where $X$ is the given a point cloud and $Z$ is the reference metric space.

- Select a covering $\mathcal{U}$ of $Z$.

- If $\mathcal{U} = \{U_\alpha\}_{\alpha \in A}$, then construct the subsets $X_\alpha = f^{-1}(U_\alpha)$.

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Representing Shape

## The Statistical version of Mapper

- Define a reference map $f : X \to Z$, where $X$ is the given a point cloud and $Z$ is the reference metric space.

- Select a covering $\mathcal{U}$ of $Z$.

- If $\mathcal{U} = \{U_\alpha\}_{\alpha \in A}$, then construct the subsets $X_\alpha = f^{-1}(U_\alpha)$.

- The analog of taking connected components in the point cloud world is clustering. Clusters form a covering of $X$ parametrized by pairs $(\alpha, c)$, where $\alpha \in A$ and $c$ is one of the clusters of $X_\alpha$.

A Short
Introduction
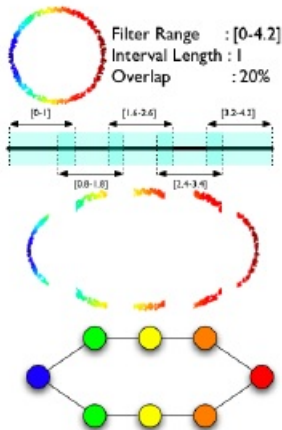to TDA
(Topological
Data Analysis)

Sara Kališnik

# Representing Shape

## The Statistical version of Mapper

- Define a reference map $f : X \to Z$, where $X$ is the given a point cloud and $Z$ is the reference metric space.

- Select a covering $\mathcal{U}$ of $Z$.

- If $\mathcal{U} = \{U_\alpha\}_{\alpha \in A}$, then construct the subsets $X_\alpha = f^{-1}(U_\alpha)$.

- The analog of taking connected components in the point cloud world is clustering. Clusters form a covering of $X$ parametrized by pairs $(\alpha, c)$, where $\alpha \in A$ and $c$ is one of the clusters of $X_\alpha$.

- Construct a graph whose vertex set is the set of all possible such pairs $(\alpha, c)$, and where an edge connects $(\alpha_1, c_1)$ and $(\alpha_2, c_2)$ if and only if the corresponding clusters have a point in common.

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Representing Shape

## The Statistical version of Mapper



Example:
Consider point cloud data which is sampled from a noisy circle in $\mathbb{R}^2$, and the filter $f(x) = ||x - p||^2$, where $p$ is the left most point in the data.

Vertices are colored by the average filter value.

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Mapper

### The Miller-Reaven diabetes study

G.M. Reaven and R.G. Miller conducted a diabetes study at Stanford in the 1970'.

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Mapper

### The Miller-Reaven diabetes study

G.M. Reaven and R.G. Miller conducted a diabetes study at Stanford in the 1970'.

145 patients were included and six quantities were measured: age, relative weight, fasting plasma glucose, area under the plasma glucose curve for the three hour glucose tolerance test(OGTT), area under the plasma insulin curve for OGTT, steady state plasma glucose response.

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Mapper

### The Miller-Reaven diabetes study

G.M. Reaven and R.G. Miller conducted a diabetes study at Stanford in the 1970'.

145 patients were included and six quantities were measured: age, relative weight, fasting plasma glucose, area under the plasma glucose curve for the three hour glucose tolerance test(OGTT), area under the plasma insulin curve for OGTT, steady state plasma glucose response.
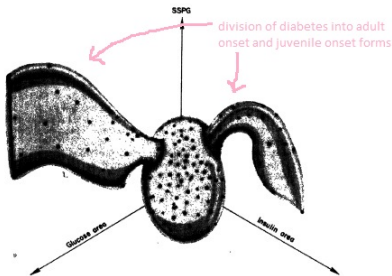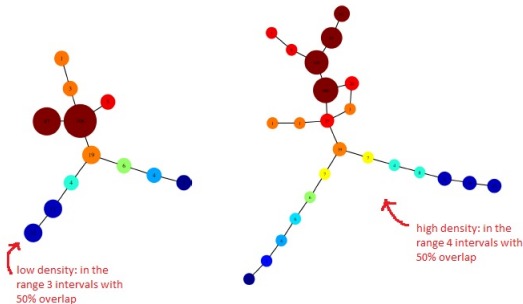
A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Mapper

### The Miller-Reaven diabetes study

If we take the filter to be a density estimator, we get the following representations for two different resolutions:



low density: in the
range 3 intervals with
50% overlap

high density: in the
range 4 intervals with
50% overlap

Red is indicative of high density, and blue of low. The size of the node and the number indicate the size of the cluster.

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Mapper

### The Miller-Reaven diabetes study

If we take the filter to be a density estimator, we get the following representations for two different resolutions:
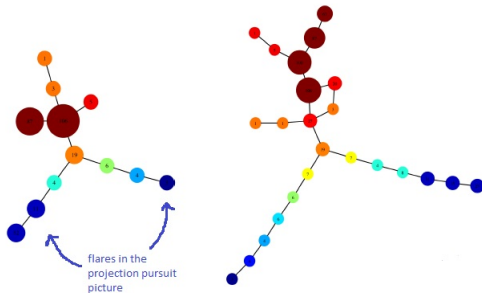


flares in the projection pursuit picture

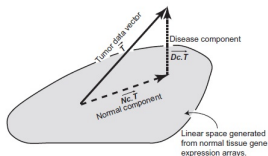Red is indicative of high density, and blue of low. The size of the node and the number indicate the size of the cluster.

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

Mapper

Breast cancer data
What should the filter be?

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Mapper

Breast cancer data
What should the filter be?



- Take linear combinations of normal expression data and denote the subspace they span by $\mathcal{N}$.

- Decompose the original data - vector $\vec{T}$ into normal-like expression, $N\vec{c}.T$, which is the projection onto $\mathcal{N}$.

- The disease, deviation $D\vec{c}.T$ from normal-like expression, is defined to be the difference between diseased tissue expression and normal-like expression.

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

Mapper

### Breast cancer data
The family of functions we take as filters is

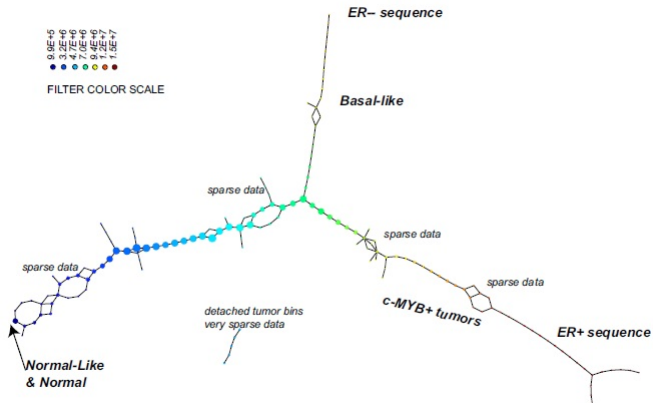$$f_{p,k}(\vec{V}) = [\sum |g_r|^p]^{\frac{k}{p}}$$

where $\vec{V} = \langle g_1, g_2, \ldots, g_s \rangle$ and coordinates $g_i$ are individual genes.

If $k = 1$, $p = 2$, the function computes standard (Euclidean) norm of a vector.

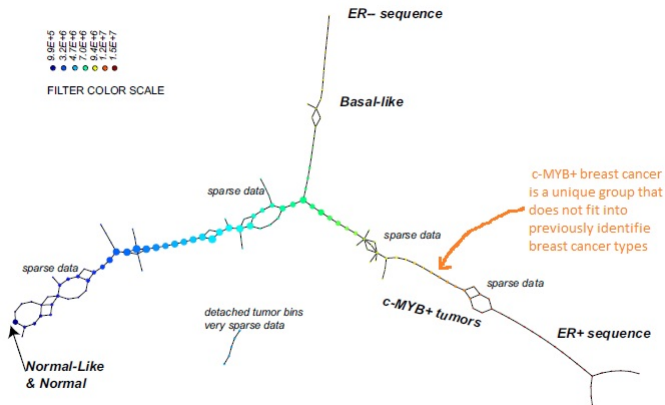Essentially, all these different filter functions, $f_{p,k}$, measure the overall amount of deviation from the normal state.

The effect of the different choices of $p$ determining the choice of $L^p$ norm is that, for larger values of $p$ the weight of genes with larger expression levels is greater.

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Mapper

## Breast cancer data

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Mapper

## Breast cancer data



FILTER COLOR SCALE

ER– sequence

Basal-like

sparse data

sparse data

sparse data

c-MYB+ breast cancer
is a unique group that
does not fit into
previously identifie
breast cancer types

c-MYB+ tumors

ER+ sequence

detached tumor bins
very sparse data

sparse data

Normal-Like
& Normal

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Mapper

## Breast cancer data



FILTER COLOR SCALE

ER– sequence

Basal-like

c-MYB+ breast cancer
is a unique group that
does not fit into
previously identifie
breast cancer types

sparse data

sparse data

sparse data

sparse data

detached tumor bins
very sparse data

c-MYB+ tumors

ER+ sequence

Normal-Like
& Normal

Both ER+ tumors (Estrogen Receptor positive) showed a 100%

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Mapper

## Clustering versus Mapper

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Mapper

## Clustering versus Mapper



c-MYB+ groups scattered across
several clusters

A Short
Introduction
to TDA
(Topological
Data Analysis)

Sara Kališnik

# Representing Shape

### Type 2 Diabetes

Current clinical definitions classify diabetes into three major subtypes: type 1 diabetes (T1D), T2D, and maturity-onset diabetes of the young.

Differences among T2D patients suggest several T2D subtypes.

Li Li, Wei-Yi Cheng, Benjamin S. Glicksberg, Omri Gottesman, Ronald Tamler, Rong Chen, Erwin P. Bottinger, and Joel T. Dudley (Icahn School of Medicine at Mount Sinai) use a topology-based approach.
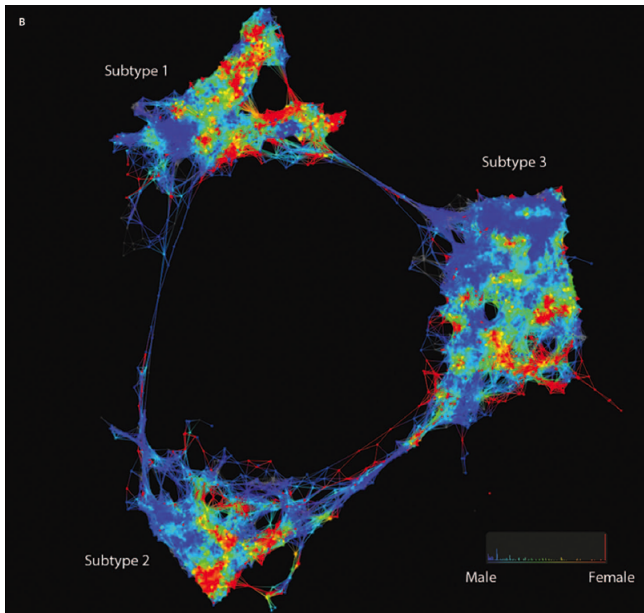
Subtype 1

Subtype 3

Subtype 2

Male                    Female

# Applied Algebraic Topology Research Network

I am one of the co-director of the Applied Algebraic Topology Research Network, which hosts a weekly Online Seminar. Recordings of our seminar are available at our YouTube Channel, which has over 6000 YouTube subscribers.



AATRN
Applied Algebraic Topology
Research Network

Applied Algebraic Topology Research Network

**Applied Algebraic Topology Network**

@aatrn1  5.95K subscribers  536 videos

This is the YouTube channel for the Applied Algebraic Topology Research …  >

aatrn.net **and 1 more link**

HOME   VIDEOS   **PLAYLISTS**   COMMUNITY   CHANNELS   ABOUT

Customize channel

Created playlists