A tutorial for
topological
data analysis
1

Otto van
Koert

Some
applications

# A tutorial for topological data analysis 1

## Otto van Koert

Seoul National University, South Korea

- Topological data analysis (TDA) is a relatively new field in mathematics with the promise of many potential new applications.

- Philosophy summarized by Carlsson as: data has shape, and shape has meaning.

- Today: Mapper and persistent homology, and some sample applications.

A tutorial for
topological
data analysis
1

Otto van
Koert

Some
applications

# Main philosophy

- Geometry and topology are subfields of mathematics that are concerned with shape

- Geometry measures quantitative properties: distance, volume, etc

- Topology is concerned with more qualitative properties, and algebraic topology transforms these into computable algebraic language.

- Applications of these techniques should be backed up with the appropriate statistics. Today we ignore these statistical issues.

A tutorial for
topological
data analysis
1

Otto van
Koert

Some
applications

# Shape of data

Data has sometimes a rather simple shape as indicated below.
Regression is an excellent technique in such cases.



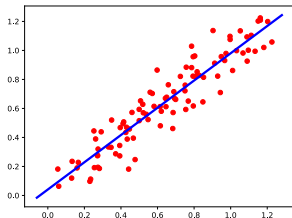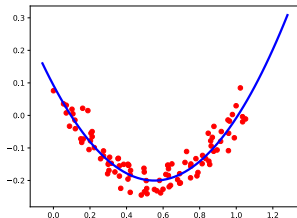Figure: Noisy linear data

Figure: Noisy quadratic data

A tutorial for
topological
data analysis
1

Otto van
Koert

Some
applications

# Other shapes

More complicated shapes appear, too.



Figure: the (transformed) amount versus time of the transaction

This shape appears in transaction data: the interpretation is obvious: each of the seven clusters corresponds to a week day: the smaller clusters are Saturday and Sunday.

A tutorial for
topological
data analysis
1

Otto van
Koert

Some
applications

# Clustering

Statistics is well-equipped to deal with the previous example: the subfield of clustering deals with this particular problem. For example, most of the following methods will work

- hierarchical clustering (single-linkage, complete linkage)
- $k$-means
- distribution-based methods such as expectation-maximization (EM)
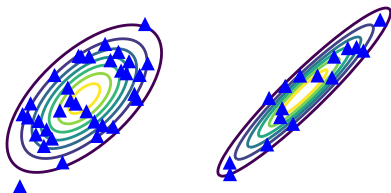- density-based methods (DB-scan)



Figure: Clusters found with EM

A tutorial for
topological
data analysis
1

Otto van
Koert

Some
applications

# More complicated shapes



Figure: Several of the above clustering algorithms will fail here



Figure: (fake) transaction data: many categorical variables

A tutorial for
topological
data analysis
1

Otto van
Koert

Some
applications

The previous example highlights some of the problems with modern data.

- a large number of data points
- many different coordinates (or features) with no direct meaning.
- related: distance does not have a direct meaning.

The last two points mean that direct geometric methods may give weak or unstable results. TDA complements these methods, because

- topology is by definition independent of the metric and
- shapes have different behavior at different scales: this will be captured with the notion of *persistence*, which relies on *functoriality*.

# Topology

We will need a couple of notions of topology, such as connectedness (previous example). We will also need a suitable class of topological spaces with a good decomposition

1. CW-complexes: preferred in algebraic topology. Examples include general graphs
2. Simplicial complexes: decomposes a topological space into many, simple pieces. This class also includes many graphs, but not all.

Simplicial complexes seem to be most useful for computational work; we will use that class after a quick review.

A tutorial for
topological
data analysis
1

Otto van
Koert

Some
applications

# Simplex as a generalized triangle

### Definition

A geometric *k*-**dimensional simplex** is the topological space

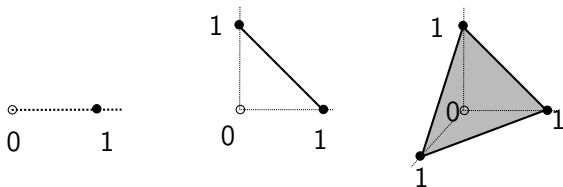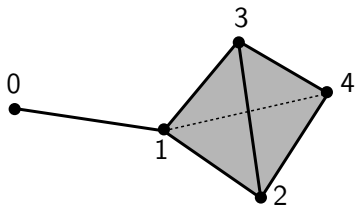$$\Delta_k := \{(x_0, \ldots, x_k) \in \mathbb{R}^{k+1} \mid \sum_{i=0}^{k} x_i = 1, \ x_i \geq 0\}$$



Figure: A 0-simplex, a 1-simplex and a 2-simplex:
a point, a line segment and a triangle

A tutorial for
topological
data analysis
1

Otto van
Koert

Some
applications

# Simplicial complexes and abstract
# simplicial complexes

- Roughly put, a simplicial complex is a topological space constructed with simplices using combinatorial "gluing recipe". See the figure below.

- An abstract simplicial complex is this combinatorial recipe; it is ideal for computers, but actually does not need any geometry.

A tutorial for
topological
data analysis
1

Otto van
Koert

Some
applications

# Simplicial complexes: Lego-like construction descriptions of spaces

### Definition (Slightly confusing)

An **abstract simplicial complex** is a (finite) collection of (finite) sets $X$ such that if $x \in X$, and $y \subset x$, then $y \in X$.

### Example

$X = \{\{0\}, \{1\}, \{2\}, \{0, 1\}\}$
We see that the definition holds.

For practical purposes, it is useful to think of
$X = X_0 \cup X_1 \cup \ldots \cup X_n$, where

- $X_k$ consists of sets with $k + 1$ elements: $X_k$ parametrizes the $k$-simplices.

In the example we have

$$X_0 = \{\{0\}, \{1\}, \{2\}\}, \quad X_1 = \{\{0, 1\}\}$$

A tutorial for
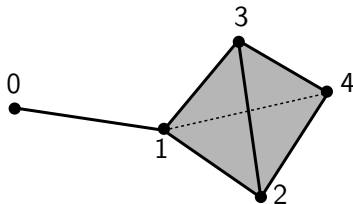topological
data analysis
1

Otto van
Koert

Some
applications

Here is a more complicated example of an abstract simplicial complex.

$$X_0 = \{\{0\}, \{1\}, \{2\}, \{3\}, \{4\}\}$$

$$X_1 = \{\{0,1\}, \{1,2\}, \{1,3\}, \{1,4\}, \{2,3\}, \{2,4\}, \{3,4\}\}$$

$$X_2 = \{\{1,2,3\}, \{1,2,4\}, \{1,3,4\}, \{2,3,4\}\}$$

We can construct a topological space out of this by replacing each $k$-simplex by a geometric $k$-simplex.



### Remark
*For later purposes (orientations) we will always order vertices in a simplex by increasing index.*

A tutorial for
topological
data analysis
1

Otto van
Koert

Some
applications

# Geometric realization:

Given an abstract simplicial complex $X = \{X_n\}_{n=0}^N$, we will define the **geometric realization** of $X$. Intuitively, this is just the shape built from the Lego description.
Here is a simple way to make it explicit:

1. order $X_0$

2. choose $N$ sufficiently large and an embedding $i : X_0 \to \mathbb{R}^N$ such that $\{i(x)\}_{x \in X_0}$ are linearly independent (can be weakened)

3. for each $\sigma \in X_k$ we get a map $f_\sigma : \Delta_k \to \mathbb{R}^N$ by linear combination. Namely if $x \in \Delta_k$, then $x = \sum_{j=0}^k t_j e_j$, where $e_j$ is the standard basis of $\mathbb{R}^{k+1}$.
   Put $f_\sigma(x) = \sum_j t_j i(\sigma[j])$, where $\sigma[j]$ is the $j$-th point of $\sigma$.

4. the geometric realization is $|X| = \cup_{\sigma \in X} f_\sigma(\Delta_{|\sigma|})$

A tutorial for
topological
data analysis
1

Otto van
Koert

Some
applications

### Remark
*In step 2 we asked for linear independence to prevent unwanted intersections between the images of different simplices. In most cases, linear independence is much stronger than necessary.*

### Remark
*We will never really need geometric realization. The abstract simplicial complex suffices for all computational work.*

### Remark
*Ordering the vertices, i.e. $X_0$, is the standard way to deal with orientations (which we only discuss implicitly).*

### Definition
A (finite) simplicial complex is a topological space that is homeomorphic to the geometric realization of a (finite) abstract simplicial complex.

A tutorial for
topological
data analysis
1

Otto van
Koert

Some
applications

# Simplicial complexes and
# compression

Many interesting (but not all) topological spaces admit a
simplicial structure. If this is the case, we can see such a
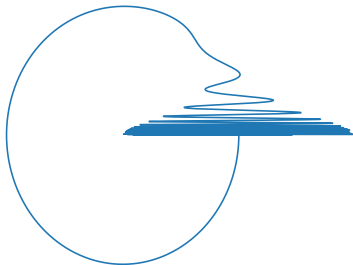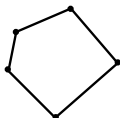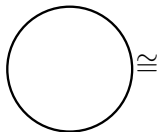simplicial structure as a compressed representation:



Figure: Circle as a simplicial
complex



Figure: The quasi-circle has no
simplicial structure

A tutorial for
topological
data analysis
1

Otto van
Koert

Some
applications

# Nerve of a covering

Let $\mathcal{U} = \{U_i\}_{i \in I}$ be a covering of a topological space $X$.

## Definition

The **nerve of the covering** $\mathcal{U}$ is the simplicial complex $N(\mathcal{U})$, whose simplices are defined as follows:

- vertices are $N(\mathcal{U})^0 = \{u_i\}_{i \in I}$, so $U_i \in \mathcal{U}$ gives one vertex.
- $[u_0, \ldots, u_k]$ forms a $k$-simplex when $\cap_{j=0}^{k} U_j \neq \emptyset$.
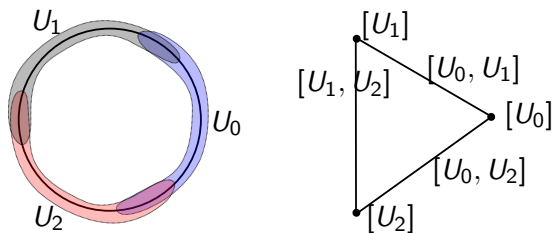


Figure: A covering of the circle and its nerve

A tutorial for
topological
data analysis
1

Otto van
Koert

Some
applications

# Nerve theorem

## Theorem (Nerve theorem)

*Suppose that $X$ is a paracompact topological space, and assume that $\mathcal{U} = \{U_i\}_i$ is a good cover, meaning that any finite intersection*

$$U_{i_0} \cap \ldots \cap U_{i_k}$$

*is either empty or contractible. Then the nerve $\mathcal{N}(\mathcal{U})$ is homotopy equivalent to $X$.*

Roughly speaking, this theorem says that the nerve of a good cover can be deformed to the original space.

## Remark

*The mathematical-philosophical background (more later) is the following. Most invariants from algebraic topology depend only on the homotopy type of a space. So replacing a space by a homotopy equivalent one that is easier, helps computations.*

A tutorial for
topological
data analysis
1

Otto van
Koert

Some
applications

# A brief excursion to the Reeb graph

Suppose that $X$ is a topological space, and $f : X \to \mathbb{R}$ a continuous function. Define an equivalence relation on $X$ by

$x \sim_f x'$ if and only if $x, x'$ are in the same component of $f^{-1}(y)$.

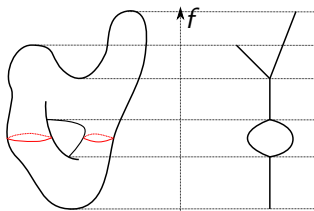Define the **Reeb graph** of $(X, f)$ as

$$R_f(X) := X/ \sim_f$$



Figure: Reeb graph of a function

A tutorial for
topological
data analysis
1

Otto van
Koert

Some
applications

# Pullback covers

We have the following relation between the Reeb graph and the nerve construction. Fix a continuous function $f : X \to \mathbb{R}$. Call $f$ a filter function or lens.

- Cover $f(X)$ by intervals $I_i$.
- The sets $U_i = f^{-1}(I_i)$ form a cover (almost never good)
- To improve our "chances", we decmpose $U_i$ into its connected components $U_i = \cup_j C_{i,j}$. If the space $X$ is reasonable (locally connected), then $C_{i,j}$ are open.

This results in an open cover $X = \cup_{i,j} C_{i,j}$.

- we can now look at the nerve construction for this cover
- and to the Reeb graph of $f$

If the cover by intervals $I_i$ is sufficiently fine, we get the "same" result.

A tutorial for
topological
data analysis
1

Otto van
Koert

Some
applications

# Reeb graph and mapper

### Remark

*The Reeb graph captures a rough version of the topology: a lot of information is lost, but some is retained in a graph, which is computationally easier to deal with.*

We cannot directly apply the Reeb graph to point cloud data; with only finitely many points in the cloud, most level sets are empty. Hence replace a level-set by the preimage of an interval.
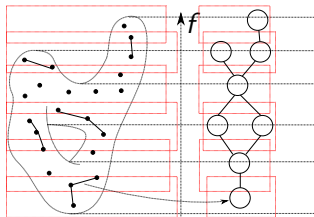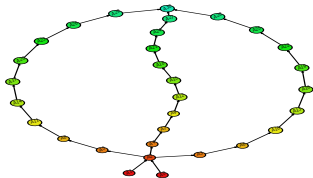


Figure: An analog of the Reeb graph

A tutorial for
topological
data analysis
1

Otto van
Koert

Some
applications

# ..Mapper

1. fix a data set and a function (filter function) from the data set to $\mathbb{R}$.
2. cover $\mathbb{R}$ with intervals that overlap on a smaller interval (gain)
3. put points whose filter value lies in interval $I_b$ in bin $b$.
4. cluster points in each bin $b$: these clusters are the vertices of the mapper
5. connect the mapper vertices if the corresponding clusters share a point.

A tutorial for
topological
data analysis
1

Otto van
Koert

Some
applications

# Multi-dimensional mapper

For the higher-dimensional version of mapper, the nerve construction is our guide. Suppose we want to construct a $k$-dimensional mapper for a pointcloud $X$. Then

1. choose filter functions $f_i : X \to \mathbb{R}$ for $i = 1, \ldots k$. Collect them in $F : X \to \mathbb{R}^k$

2. cover $F(X)$ with overlapping boxes $U_i$

3. to obtain the analog of the pullback cover, cluster data points in each $U_i$. Call the clusters $C_{i,j}$

4. The clusters form the vertices of mapper

5. There is an edge if two clusters $C_{i,j}$ and $C_{i',j'}$ overlap (i.e. share a data point) This can be generalized to several clusters to get higher-dimensional simplices.

A tutorial for
topological
data analysis
1

Otto van
Koert

Some
applications

The main purposes of mapper are:

- visualization
- exploratory data analysis

It is flexible and can deal with large and diverse data sets: mixed numerical and categorical data can be dealt with. One only needs a dissimilarity measure.
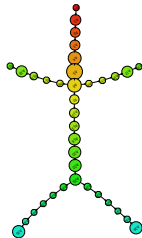


Figure: A model of a human



Figure: Mapper of the first PCA-function for a human model; PCA stands for principal component analysis

A tutorial for
topological
data analysis
1

Otto van
Koert

Some
applications

# Choices

There are lots of choices in mapper. This gives it a lot of flexibility, but this comes at a price, namely complexity.

- We need to choose filter functions of interest
- The number of bins, and the epsilon parameters need to be chosen.

Here it should be pointed out that statistics gives us some guidance.

- more bins speeds up the mapper algorithm
- but more bins results in fewer points per bin, which makes the statistical data (mean, etc per bin) less reliable.
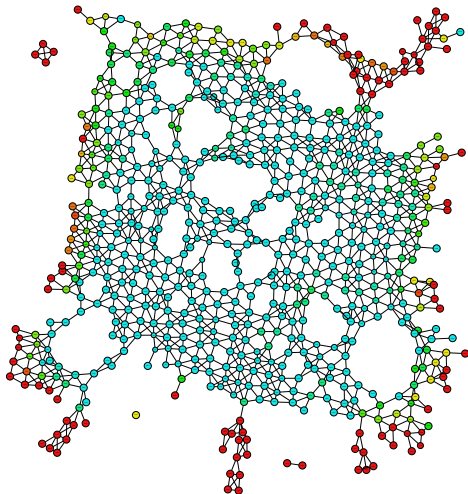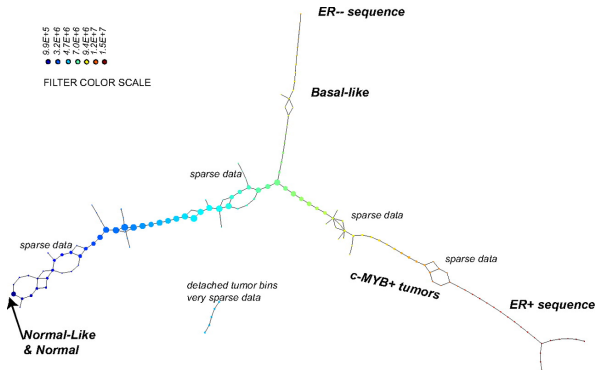
A tutorial for
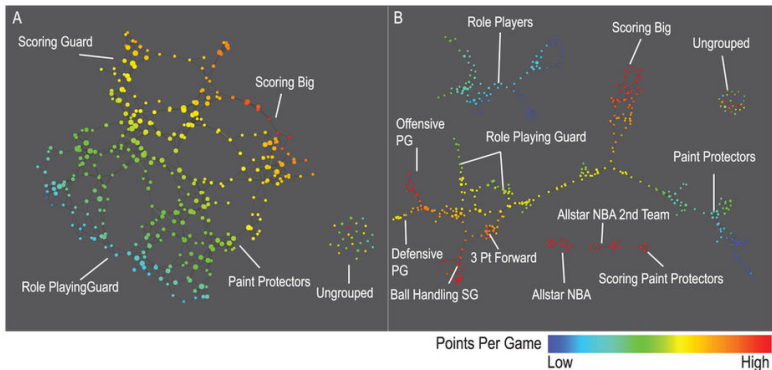topological
data analysis
1

Otto van
Koert

Some
applications

# Some applications



Figure: Mapper of medical data with two filter functions, and colored by a third filter function (disease: yes or no)

A tutorial for
topological
data analysis
1

Otto van
Koert

Some
applications

The following picture is taken from
*Topology based data analysis identifies a subgroup of breast cancers with a unique mutational profile and excellent survival* by Monica Nicolau, Arnold J. Levine, and Gunnar Carlsson, PNAS April 26, 2011. 108 (17) 7265-7270;

The following picture is taken from *Extracting insights from the shape of complex data using topology* by P. Y. Lum, G. Singh, A. Lehman, T. Ishkanov, M. Vejdemo-Johansson, M. Alagappan, J. Carlsson and G. Carlsson, Nature, Scientific Reports volume 3, Article number: 1236

A tutorial for
topological
data analysis
1

Otto van
Koert

Some
applications

Thank you
감사 합니다.