

# SCHEMA INFERENCE

## ALGEBRAIC TOPOLOGY FOR SCHEMA DISCOVERY

LUCIANO MELODIA

CHAIR OF COMPUTER SCIENCE 6  
FRIEDRICH-ALEXANDER UNIVERSITY

[LUCIANO.MELODIA@FAU.DE](mailto:LUCIANO.MELODIA@FAU.DE)

JULY 17, 2019



# OVERVIEW

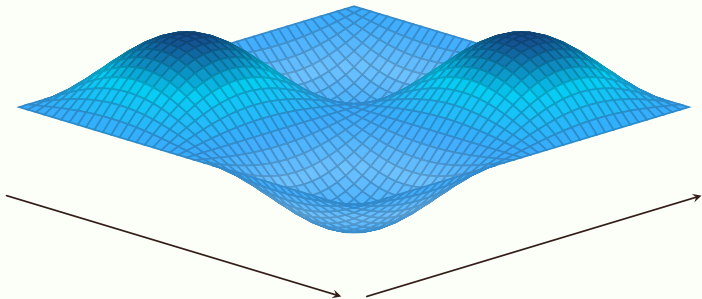
WHAT WILL I PRESENT TODAY?

- Motivation
- Simplicial Complexes
- Persistent Homology
- Current Projects
- Future Work (joint with Siemens on real data)

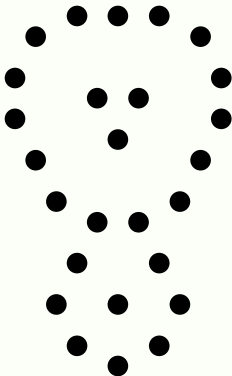
# **PART I: MOTIVATION**

# WHERE DOES DATA COME FROM?

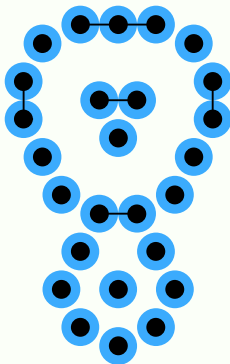
WE ARE INTERESTED IN THE SHAPES OF MANIFOLDS.



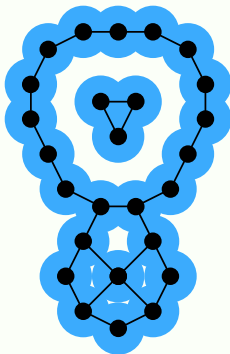
What is the 'shape' of the world where our data comes from?



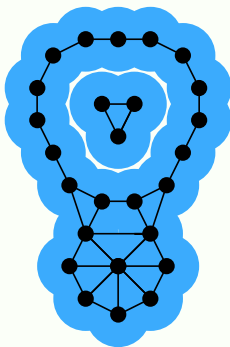
This point cloud has no form in its own right in the technical sense, but we clearly see two 'rings' connected to each other. In the smaller lower 'ring' there is a point, in the larger upper ring a 'triangle'.



Imagine every point as a balloon. We now inflate this balloon by one  $r = 0.2$  and more distinct structures become visible in space.

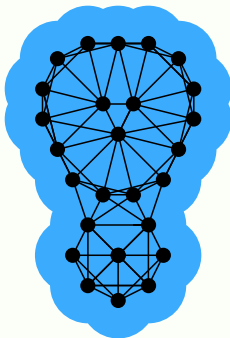


As soon as two balloons touch, we create an edge between the two center points. This  $r = 0.4$ .



$$r = 0.6.$$





This is the 1-skeleton with an  $r = 0.8$ .

## WHAT DID WE JUST SEE?

The formation of a 1-dimensional graph by successively increasing the radius of the balls has revealed the underlying structure.

How can this be formulated more precisely?

This is the area of mathematics which deals with invariant properties of high dimensional geometric objects or spaces.

Simple invariants:

- Dimension:  $\mathbb{R}^2 \neq \mathbb{R}^3$  because  $2 \neq 3$ .
- Determinant: is  $\det(A) = \det(B)$  for two matrices  $A, B$ , then they are considered to be similar.

We believe that different data comes from different spaces. How can we detect this and sort the data accordingly?

We don't want to distinguish data only by their 'holes', therefore we would like to introduce a 'magnifying glass' which shall allow us to capture the **structure of a space** in different granularity.

Which structures do we observe?

## **PART II: SIMPLICIAL STRUCTURES**

# REAL-WORLD MULTIVARIATE DATA

TIME SERIES, NON-ANALYTIC OR ANISOTROPIC DATA, NOISE AND MORE ...

- Often unstructured point clouds.
- $n$  items with  $d$  attributes gives  $n \times d$  matrices.
- Non-random often anisotropic sampling from  $\mathbb{R}^d$ .

## Manifold assumption.

There is a  $n$ -dimensional Riemannian manifold  $M$  with  $n \ll d$  where our data comes from.

# SIMPLEX

OUR STRUCTURE TO WORK WITH.

Given a set  $X = \{x_0, \dots, x_k\} \subset \mathbb{R}^d$  of  $k + 1$  affinely independent points, the  $k$ -dimensional simplex  $\Delta$  spanned by  $X$  is the set of convex combinations

$$\sum_{i=0}^k \lambda_i x_i, \quad \text{with} \quad \sum_{i=0}^k \lambda_i = 1 \quad \text{and} \quad \lambda_i \geq 0. \quad (1)$$

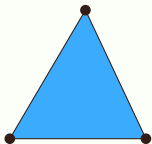




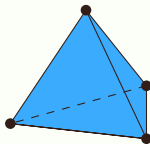
$\Delta^0$  or 0-simplex



$\Delta^1$  or 1-simplex



$\Delta^2$  or 2-simplex



$\Delta^3$  or 3-simplex

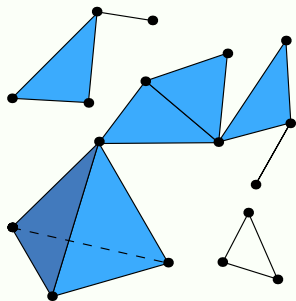
# SIMPLICIAL COMPLEX

GLUING TOGETHER PIECES OF INCREASING DIMENSION.

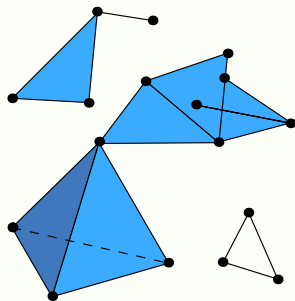
- Gluing together pieces of increasing dimensions.
- Preserve the dimension of the simplex.
- The dimension of the simplicial complex is the dimension of the  $n$ -simplex.

A *simplicial complex*  $K$  is a topological space realized as a union of any collection of simplices  $\Delta$  with the following properties:

- Any face of a simplex  $\Delta$  is also in  $\Delta$ .
- The intersection of any two simplices of  $\Delta$  is also in  $\Delta$ .



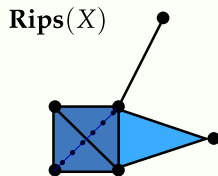
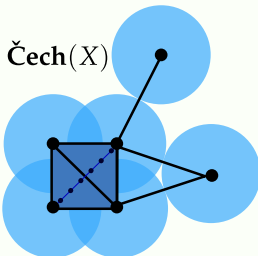
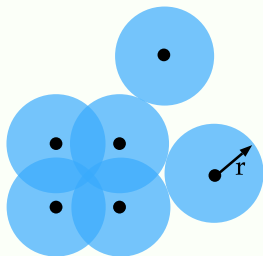
A valid simplicial complex.



Invalid construction.

# CONVERTING DATA TO SIMPLICIAL COMPLEXES

I.E. ČECH-, VIETORIS-RIPS- OR WITNESS-COMPLEX.



Creating a Čech-complex: Use any distance measure  $dist(\cdot, \cdot)$  on your manifold and determine a threshold  $r$ . Create an edge between two points if  $dist(p_i, p_j) \leq r$ .

# FINDING PATTERNS

## BOUNDARIES.

- Boundary is easy to describe.
- Geometrically obvious.

The boundary of a line segment are the two endpoints.

The boundary of a triangle is the union of its edges.

The boundary of a tetrahedron is the union of triangular faces.

⋮

# FINDING PATTERNS

## BOUNDARIES.

- A boundary itself has no boundary.
- Being boundariless in finite dimensions coincides with the intuition of a loop.
- We define these things in algebraic terms.

## PART III: CALCULATING $H_k(X)$

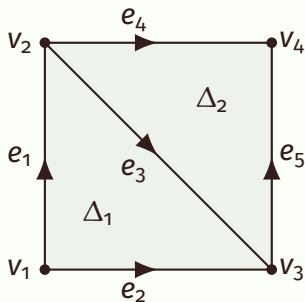
# REALIZE SIMPLICES AS ALGEBRAIC OBJECTS

$X_k$  is the set of  $k$ -simplices of the simplicial complex  $X$ .  
The *chain group*  $C_k(X)$  is the  $\mathbb{Q}$ -vector space with  $X_k$  as a basis.

## Chain groups

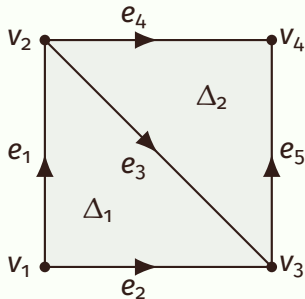
Let  $\Delta^k$  and  $\Delta^{k'}$  be two  $k$ -simplices. Then we define a bunch of new 'chains' as all possible sums and scalar multiples of the simplices, i.e.  $\lambda_1 \Delta^k + \lambda_2 \Delta^{k'}$ . This is denoted in general as  $\text{span}(\Delta_1^k, \dots, \Delta_n^k)$ .





$C_0(X)$  is the span of the vertices  $\text{span}(v_1, v_2, v_3, v_4)$ . Having two copies of  $v_2$  and four of  $v_3$  yields this algebraic equation  $2v_2 + 4v_3$ .

There are meaningless equations such as  $\frac{1}{2}v_4 - v_2 - \frac{3}{4}v_5$  which exist for their own right.



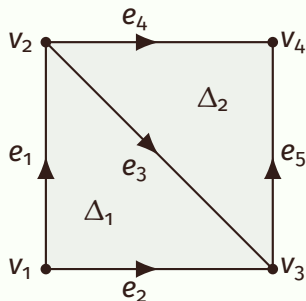
$C_1(X)$  is therefore the span of the set  $\{e_1, e_2, e_3, e_4, e_5\}$ . We can talk about a path, as the edges are line segments.

Again some realizations are geometrically ridiculous:  $-e_5 + e_4 + e_1 + e_3$ .

- The chain group  $C_2(X) = \text{span}(\Delta_1, \Delta_2)$ .
- What is a path of triangles?
- Geometric analogies become tenuous in high dimensions.
- Algebra helps. Find low dimensional intuitions and generalize them.

- Chain groups are linear groups.
- Their elements are vector spaces.
- They induce linear maps  $\partial : C_k(X) \rightarrow C_{k-1}(X)$ .
- Single vertices are boundariless, i.e.  $\partial(v) = 0$ .

# THE BOUNDARY OPERATOR



$$\partial(v_1, v_2) = v_2 - v_1$$

$$\begin{aligned}\partial(e_1 + e_4 - e_5 - e_3) &= \partial e_1 + \partial e_4 - \partial e_5 - \partial e_3 = \\ &= (v_2 - v_1) + (v_4 - v_2) - (v_4 - v_3) - (v_3 - v_2) = \\ &= v_2 - v_1.\end{aligned}$$

# THE BOUNDARY OPERATOR

## RESULTS.

- The endpoint of  $v_2$  minus the starting point of  $v_1$ . Each successive edge cancels out the ending vertex of the edge before it.
- Each starting edge cancels the starting edge after it.
- In case of a loop the boundary operator becomes 0.
- The boundary operator is given by an alternating sum.

# THE BOUNDARY OPERATOR

GENERALIZED FORM.

The *boundary operator* or sometimes *boundary map*

$\partial_k : C_k(X) \rightarrow C_{k-1}(X)$  is given by:

$$\partial_k(\Delta^k) = \sum_{i=0}^k (-1)^i \Delta^k \Big|_{[v_0, \dots, v_{i-1}, \hat{v}_i, v_{i+1}, \dots, v_k]} \quad (2)$$

- Boundary operators can be composed:  $\partial \circ \partial = \partial^2$ .
- Every chain which is a boundary of higher-dimensional chains is boundariless, i.e.  $\partial^2 \mathbb{Z} = 0$ .
- Low dimensional examples make sense and can be extended.
- Boundariless 'cycles' correspond to loops.



# OBSERVATIONS

... IN THE LANGUAGE OF ALGEBRAIC TOPOLOGY.

- The kernel is defined as

$$\ker \partial_k := \partial_k^{-1}(e_{C_{k-1}(X)}) = \{z \in C_k(X) : \partial_k(z) = e_{C_{k-1}(X)}\}.$$

- A  $k$ -cycle is then defined as  $Z_k(X) = \ker(\partial_k) \subset C_k(X)$ .

- The  $k$ -boundaries  $B_k(X) = \text{im}(\partial_{k+1}) \subset Z_k(X)$ .

$$\dots \xrightarrow{\partial_{k+1}} C_k \xrightarrow{\partial_k} C_{k-1} \xrightarrow{\partial_{k-1}} \dots \xrightarrow{\partial_2} C_1 \xrightarrow{\partial_1} C_0 \xrightarrow{\partial_0} 0 \xrightarrow{\partial_{k+1}} \dots \quad (3)$$

# SIMPLICIAL HOMOLOGY GROUP

The  $k^{\text{th}}$  homology group  $H_k$  is a quotient group, defined by 'removing' cycles that are boundaries from a higher dimension:

$$H_k = Z_k/B_k = \ker \partial_k / \text{im } \partial_{k+1} \quad (4)$$

This definition gives rise to the  $k^{\text{th}}$  Betti-number:

$$\beta_k = \text{rank } H_k \quad (5)$$

Intuitively: Calculate all boundaries, remove the boundaries that come from higher-dimensional objects, and count what is left.

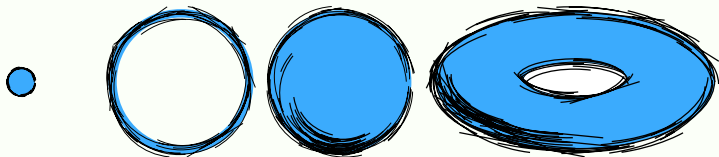
# BETTI NUMBERS

ANOTHER TOPOLOGICAL INVARIANT.

Betti numbers 'count' the number of holes in the dimensions that occur in space.

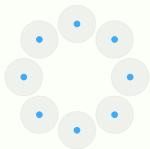
$\beta_0$  Connected components  
 $\beta_1$  Tunnels  
 $\beta_2$  Voids  
 $\vdots$   $\vdots$

Space	$\beta_0$	$\beta_1$	$\beta_2$
Point	1	0	0
Circle	1	1	0
Sphere	1	0	1
Torus	1	2	1

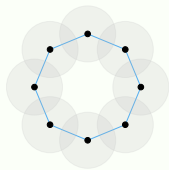


# BETTI NUMBERS

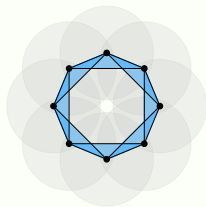
ANOTHER TOPOLOGICAL INVARIANT.



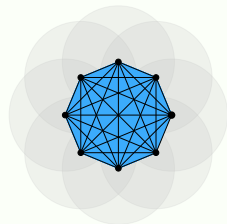
**(a)**  $r = 0.2$ .



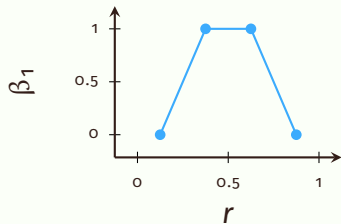
**(b)**  $r = 0.4$ .



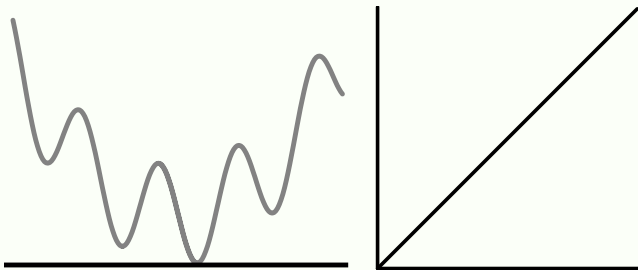
**(c)**  $r = 0.6$ .

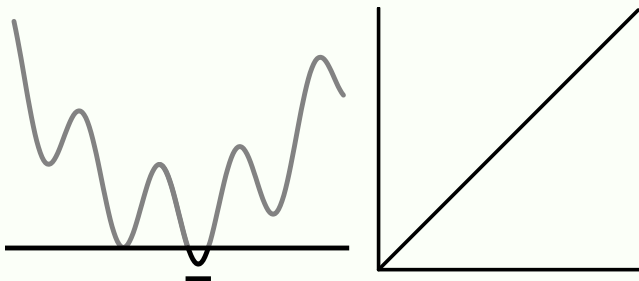


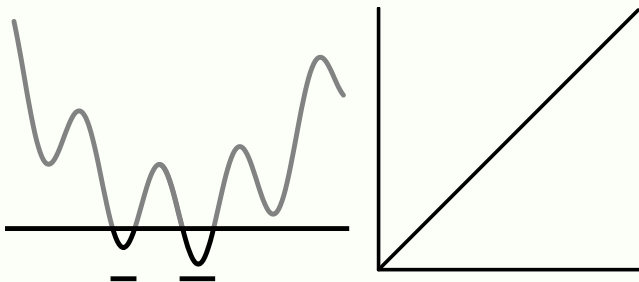
**(d)**  $r = 0.8$ .



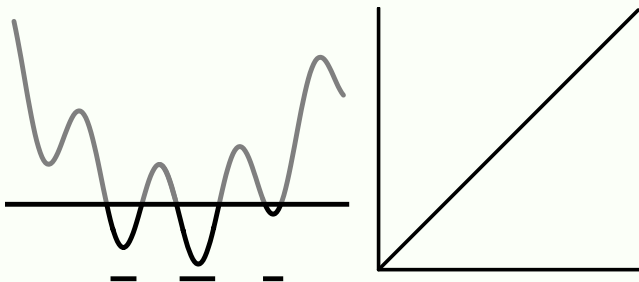
- Calculates the  $k$ -th Betti numbers.
- Gives an overview of the evolving topology of a filtration.
- Robust against noise.
- Detects dense regions in data.

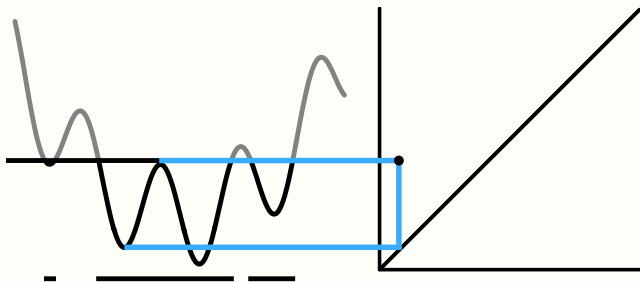


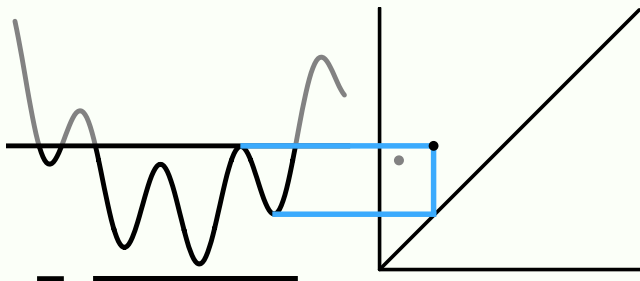


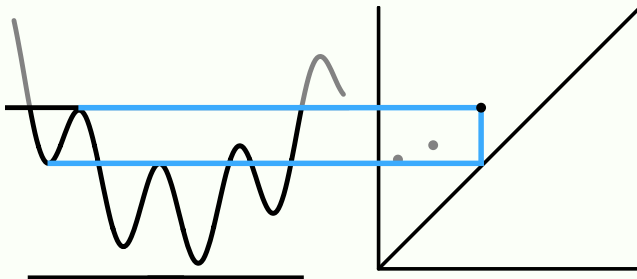


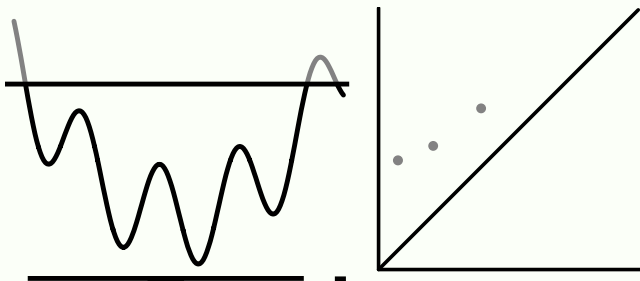


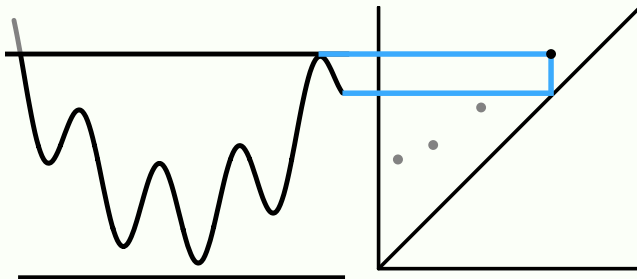


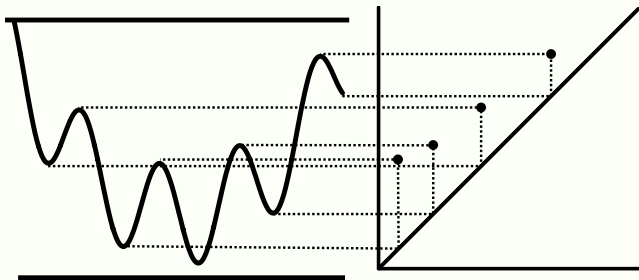






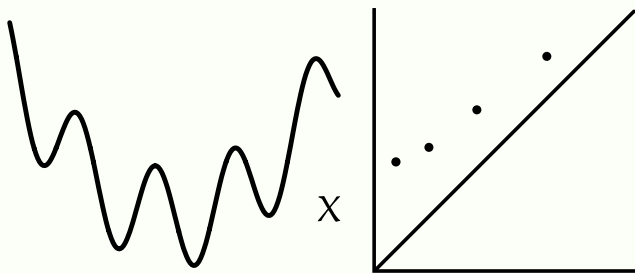






# PERSISTENT HOMOLOGY

PERSISTENCE DIAGRAM BY RIECK.

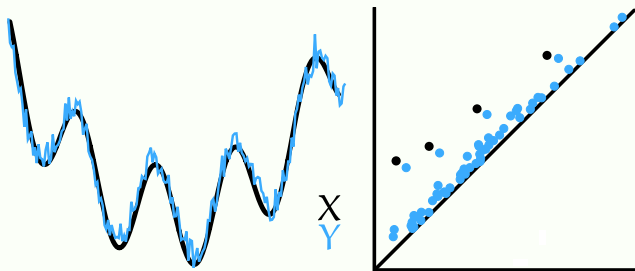


A persistence diagram of a curve.



# PERSISTENT HOMOLOGY

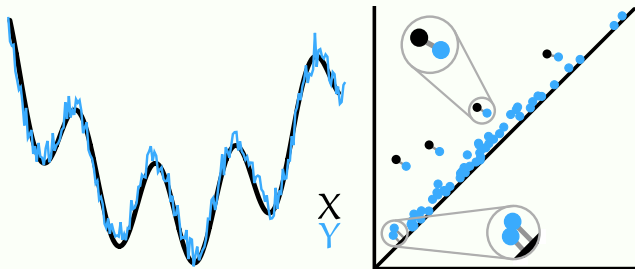
PERSISTENCE DIAGRAM BY RIECK.



Persistent homology is robust against noise in a very precise way.

# MEASURE OF DISTANCE

## BOTTLENECK DISTANCE.

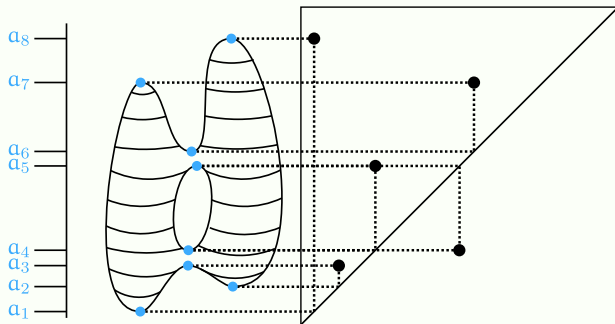


The **Bottleneck** distance compares two persistence diagrams:

$$d_B(\mathcal{P}_1, \mathcal{P}_2) = \inf_{\varphi} \sup_{x \in \mathcal{P}_1} \|x - \varphi(x)\|_{\infty}$$

# MORSE THEORY

## CONNECTION TO DIFFERENTIAL TOPOLOGY.



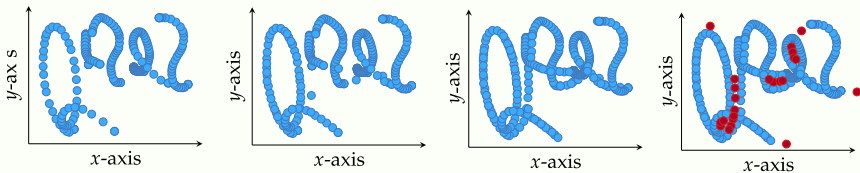
Minima in hyperplanes create new 'connected' components.  
Maxima in turn destroy them by merging connected components.  
Saddle points either create cavities or 'holes' or connect two connected components.

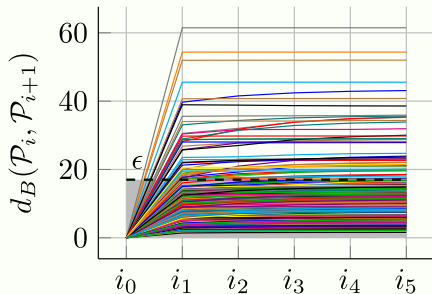
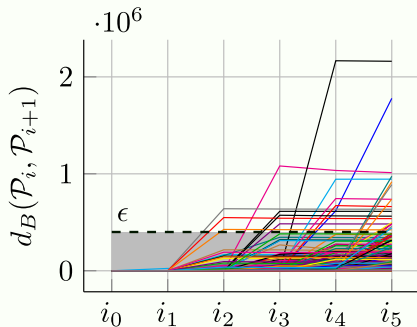
# **PART V: CURRENT PROJECTS**

# CURRENT PROJECTS

## NATURAL NEIGHBOR INTERPOLATION WITH STOP.

- Method of Natural Neighbors for interpolation.
- Topological plausible interpolation.
- Stabilizing the method with topological stopping.
- Interpolation performs until topology remains 'the same'.



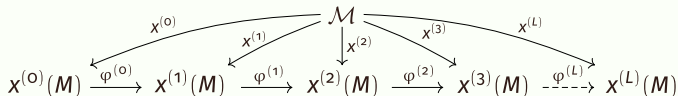


Natural Neighbor interpolation with heuristic stop  $\epsilon$ .

Different filtrations lead to impressively similar results, but are dependent on parametrization.

# CURRENT PROJECTS

## UPPER BOUND FOR RESIDUAL NEURAL NETWORK WIDTH.



### ■ **Manifold assumption:**

Data lives on low dimensional manifold.

■ What is the dimension of this manifold  $\mathcal{M}$ ?

■ How 'complex' is the topology of the manifold  $\mathcal{M}$ ?

■ How many dimension do I really need for trajectory  $\varphi^{(l)}$ ?

# CURRENT PROJECTS

SIGNATURE IDENTIFICATION JOINT WITH MICHAEL NISSEN.

- Visualization of the topology of signatures.
- Knowledge inference from persistent homology.
- Choosing appropriate filtrations.
- Interpolation of the data.
- Parametrizing neural networks with knowledge from persistent homology.
- Using residual networks to identify signatures.



## **PART IV: FUTURE WORK**

# MOTIVATION

## BEYOND MATHEMATICAL LOGIC.

Applying algebraic topology and category theory yields the category of databases as theoretical framework. Four good reasons for category theory:

- Morphisms preverse the database structure.
- Categorical constructions (limits or colimits) result in database operations (union or join).
- Restructuring and integrating data becomes easy.
- The schema has a natural geometrical structure.

Following the work of **David Spivak** a relation can be realized as simplicial complex. Now the following topics can be posed, which I'll try to elaborate:

- Application of Spivaks construction to real world data.
- Visualization of functional dependencies using TDA.
- Triangulation of relations capturing FDs.
- Clustering of large datasets according to this triangulation.

# REFERENCES



JEAN-DANIEL BOISSONNAT, FRÉDÉRIK CHAZAL, MARIETTE YVINEC  
**GEOMETRIC AND TOPOLOGICAL INFERENCE.**  
*Cambridge, 2018.*



DAVID SPIVAK  
**SIMPLICIAL DATABASES.**  
*Arxiv abs:0904.2012, 1–25, 2009.*



BASTIAN RIECK  
**PERSISTENT HOMOLOGY IN MULTIVARIATE DATA VISUALIZATION.**  
*University of Heidelberg, pages 1–307, 2017.*



BASTIAN RIECK  
**A PRIMER IN PERSISTENT HOMOLOGY.**  
*URL: <https://www.math.uni-hamburg.de/>.*



JEREMY KUN  
**HOMOLOGY THEORY - A PRIMER.**  
*URL: <https://jeremykun.com/>.*

# ACKNOWLEDGEMENT

I would like to thank **Justin Noel** and **Jan Frahm** for bringing me closer to algebraic topology and differential geometry.

Special thanks to **Bastian Rieck**, who provided me some of his graphics for this presentation with friendly permission.

Furthermore I would also like to thank my supervisor **Richard Lenz**, who breaks up complex contexts and makes me bring them into an orderly form.