

# Wasserstein convergence of Čech persistence diagrams for samplings of submanifolds

Charles Arnal, David Cohen-Steiner, Vincent Divol

November 12, 2024

# Topological data analysis and Čech persistence diagrams

Topological data analysis studies the shape of data, i.e. its geometric or topological structure.

The Čech persistence diagrams  $\text{dgm}_i(S)$  of a set  $S$  capture some aspects of its topology.

Persistence diagrams can be vectorized and used as inputs to various machine learning methods for classification or regression.

Common task in TDA: recover information regarding the shape of some set  $S$  using a finite sampling  $A \subset S$ .

**Idea:** use the persistence diagrams  $\text{dgm}_i(A)$  as estimates of  $\text{dgm}_i(S)$ .

**Problem:**  $\text{dgm}_i(A)$  is only known to converge weakly to  $\text{dgm}_i(S)$  as  $A \rightarrow S$ . In particular, machine learning-friendly features might not converge.

In our article, we show that if

- $S$  is a generic submanifold, and
- $A_n \subset S$  is a random finite sampling of cardinality  $n$ ,

then  $\text{dgm}_i(A_n) \xrightarrow{n \rightarrow \infty} \text{dgm}_i(S)$  for the stronger Wasserstein metric.

This guarantees the stability of many machine learning features.

More specifically, we show:

- A stronger version of the bottleneck theorem,
- Finiteness in expectation of the number of points in some regions of the diagram,
- Asymptotic expression for the total persistence of  $\text{dgm}_i(A_n)$ ,
- Convergence of  $\text{dgm}_i(A_n)$  to  $\text{dgm}_i(S)$  for the Wasserstein distances.