

University of Oxford



DEPARTMENT OF
STATISTICS

Topological Data Analysis of Temporal Networks

by

Dimitri Lozeve

St Catherine's College

A dissertation submitted in partial fulfilment of the degree of Master of Science in
Applied Statistics

*Department of Statistics, 24–29 St Giles,
Oxford, OX1 3LB*

September 2018

Declaration of authorship

This my own work (except where otherwise indicated).

Date

Signature

ABSTRACT

Abstract here

Acknowledgements

Thank you!

Contents

Abstract	iii
1 Introduction	1
2 Topological Data Analysis and Persistent Homology	2
2.1 Homology	2
2.2 Simplicial Complexes	2
2.3 Filtrations	3
2.4 Persistent Homology	3
2.5 Topological summaries: barcodes and persistence diagrams	3
2.6 Stability	4
3 Temporal Networks	5
3.1 Definition and basic properties	5
3.2 Network partitioning	6
3.3 Persistent homology for networks	6
3.4 Zigzag persistence	6

List of Figures

List of Tables

1 *Introduction*

2 Topological Data Analysis and Persistent Homology

2.1 HOMOLOGY

Our goal is to understand the topological structure of a metric space. For this, we can use *homology*, which consists in associating for a metric space X and a dimension i a vector space $H_i(X)$. The dimension of $H_i(X)$ will give us the number of i -dimensional components in X : the dimension of $H_0(X)$ is the number of path-connected components in X , the dimension of $H_1(X)$ is the number of holes in X , and the dimension of $H_2(X)$ is the number of voids.

Crucially, these vector spaces are robust to continuous deformation of the underlying metric space (they are *homotopy invariant*). However, computing the homology of an arbitrary metric space can be extremely difficult. It is necessary to approximate it in a structure that would be both combinatorial and topological in nature.

2.2 SIMPLICIAL COMPLEXES

In order to understand the topological structure of a metric space, we need a way to decompose it in smaller pieces which, when assembled, conserve the overall organisation of the space. For this, we use a structure called a *simplicial complex*, which is a kind of higher-dimensional generalization of graphs.

The building blocks of this representation will be *simplices*, which are simply the convex hull of an arbitrary set of points. Examples of simplices include single points, segments, triangles, and tetrahedrons (in dimensions 0, 1, 2, and 3 respectively).

Definition 2.1 (Simplex). The k -dimensional simplex $\sigma = [x_0, \dots, x_k]$ is the convex hull of the set $\{x_0, \dots, x_k\} \in \mathbb{R}^d$, where x_0, \dots, x_k are affinely independent. x_0, \dots, x_k are called the *vertices* of σ , and the simplices defined by the subsets of $\{x_0, \dots, x_k\}$ are called the *faces* of σ .

We then need a way to combine these basic building blocks meaningfully so that the resulting object can adequately reflect the topological structure of the metric space.

Definition 2.2 (Simplicial complex). A *simplicial complex* is a collection K of simplices such that:

- any face of a simplex of K is a simplex of K
- the intersection of two simplices of K is either the empty set or a common face or both.

The notion of simplicial complex is closely related to that of a hypergraph. The important distinction lies in the fact that a subset of a hyperedge is not necessarily a hyperedge itself.

Using these definitions, we can define homology on simplicial complexes.

2.3 FILTRATIONS

If we consider that a simplicial complex is a kind of “discretization” of a metric space, we realise that there must be an issue of *scale*. For our analysis to be invariant under small perturbations in the data, we need a way to find the optimal scale parameter to capture the adequate topological structure, without taking into account some small perturbations, nor ignoring some important smaller features.

The ideal solution to these problems is to consider all scales at once: this is the objective of *filtered simplicial complexes*.

Definition 2.3 (Filtration). A *filtered simplicial complex*, or simply a *filtration*, K is a sequence $(K_i)_{i \in I}$ of simplicial complexes such that:

- for any $i, j \in I$, if $i < j$ then $K_i \subseteq K_j$,
- $\bigcup_{i \in I} K_i = K$.

2.4 PERSISTENT HOMOLOGY

We can now compute the homology for each step in a filtration. This leads to the notion of *persistent homology*, which gives us all the information necessary to establish the topological structure of the metric space at multiple scales.

Definition 2.4 (Persistent homology). The p -th *persistent homology* of a simplicial complex $K = (K_i)_{i \in I}$ is the pair $(\{H_p(K_i)\}_{i \in I}, \{f_{i,j}\}_{i,j \in I, i \leq j})$, where for all $i \leq j$, $f_{i,j} : H_p(K_i) \mapsto H_p(K_j)$ is induced by the inclusion map $K_i \mapsto K_j$.

The functions $f_{i,j}$ allow us to link generators in each successive homology space in the filtration. Since each generator correspond to a topological feature (connected component, hole, void, etc, depending on the dimension p), we can determine whether it survives in the next step of the filtration. We can now determine when each feature is born and when it dies (if it dies at all). This representation will be dependent on the choice of basis for each homology space $H_p(K_i)$. However, by the Fundamental Theorem of Persistent Homology, we can choose base vectors in each homology space such that the collection of half-open intervals is well-defined and unique. This construction is called a *barcode*.

2.5 TOPOLOGICAL SUMMARIES: BARCODES AND PERSISTENCE DIAGRAMS

In order to interpret the results of the persistent homology computation, we need to compare the output for a particular data set to a suitable null model. For this, we need some kind of a similarity measure between barcodes and a way to evaluate the statistical significance of the results.

One possible approach for this is to define a space in which we can project barcodes and study their geometric properties. *Persistence diagrams* are an example of such a space.

Definition 2.5 (Persistence diagrams). A *persistence diagram* is the union of a finite multiset of points in $\bar{\mathbb{R}}^2$ with the diagonal $\Delta = \{(x, x) \mid x \in \mathbb{R}^2\}$, where every point of Δ has infinite multiplicity.

The diagonal Δ is added to facilitate comparisons between diagrams, as points near the diagonal correspond to short-lived topological feature, thus likely to be caused by small perturbations in the data.

We can now define several distances on the space of persistence diagrams.

Definition 2.6 (Wasserstein distance). The p -th Wasserstein distance between two diagrams X and Y is

$$W_p[d](X, Y) = \inf_{\phi: X \rightarrow Y} \left[\sum_{x \in X} d(x, \phi(x))^p \right]$$

for $p \in [1, \infty)$, and

$$W_\infty[d](X, Y) = \inf_{\phi: X \rightarrow Y} \sup_{x \in X} d(x, \phi(x))$$

for $p = \infty$, where d is a distance on \mathbb{R}^2 and ϕ ranges over all bijections from X to Y .

Definition 2.7 (Bottleneck distance). The *bottleneck distance* is defined as the infinite Wasserstein distance with d the uniform norm: $d_B = W_\infty[L_\infty]$.

Since the bottleneck distance is by far the most commonly used, we will focus on it in the following. It is symmetric, non-negative, and satisfies the triangle inequality. However, it is not a true distance, as it is fairly straightforward to come up with two distinct diagrams at bottleneck distance zero, even on multisets not touching the diagonal Δ .

2.6 STABILITY

3 Temporal Networks

3.1 DEFINITION AND BASIC PROPERTIES

In this section, we will introduce the notion of temporal networks or graphs. This is a complex notion, with many concurrent definitions and interpretations. First, we restate the standard definition of a non-temporal, static graph.

Definition 3.1 (Graph). A *graph* is a couple $G = (V, E)$, where V is a finite set of *nodes* (or *vertices*), and $E \subseteq V \times V$ is a set of *edges*. A *weighted graph* is defined by $G = (V, E, w)$, where $w : E \mapsto \mathbb{R}_+$ is called the *weight function*.

We also define some basic concepts that will be needed later on to build simplicial complexes on graphs.

Definition 3.2 (Clique). A *clique* is a set of nodes where each pair is connected. That is, a clique C of a graph $G = (V, E)$ is a subset of V such that $\forall i, j \in C, i \neq j \implies (i, j) \in E$. A clique is said to be *maximal* if it cannot be augmented by any node.

Temporal networks are defined in the more general framework of *multilayer networks*. However, this definition is much too general for our simple applications, and we restrict ourselves to edge-centric time-varying graphs. In this model, the set of nodes is fixed and doesn't change over time, whereas edges can appear or disappear at different timestamps.

Definition 3.3 (Temporal network). A *temporal network* (or graph) is a tuple $G = (V, E, \mathcal{T}, \rho)$, where:

- V is a finite set of nodes,
- $E \subseteq V \times V$ is a set of edges,
- \mathbb{T} is the *temporal domain* (often taken as \mathbb{N} or \mathbb{R}_+), and $\mathcal{T} \subseteq \mathbb{T}$ is the *lifetime* of the network,
- $\rho : E \times \mathcal{T} \mapsto \{0, 1\}$ is the *presence function*, which determines whether an edge is present in the network at each timestamp.

The *available dates* of an edge are the set $\mathcal{I}(e) = \{t \in \mathcal{T} : \rho(e, t) = 1\}$.

Temporal networks can also have weighted edges. In this case, it is possible to have constant weights (edges can only appear or disappear over time, and always have the same weight), or time-varying weights. In the latter case, we can set the domain of the presence function to be \mathbb{R}_+ instead of $\{0, 1\}$, where by convention a zero weight corresponds to an absent edge.

Definition 3.4 (Additive temporal network). A temporal network is said to be *additive* if for all $e \in E$ and $t \in \mathcal{T}$, if $\rho(e, t) = 1$, then $\forall t' > t, \rho(e, t') = 1$. Edges can only be added to the network, never removed.

3.2 NETWORK PARTITIONING

3.3 PERSISTENT HOMOLOGY FOR NETWORKS

We now consider the problem of applying persistent homology to network data. An undirected network is already a simplicial complex of dimension 1. However, this will not be sufficient to capture enough topological information: we need to introduce higher-dimensional simplices. The first possible method is to project the network on a metric space, thus transforming the network data into a point cloud data. For this, we need to compute the distance between each pair of nodes in the network (via shortest path distance for instance). This also requires the network to be connected.

Another usual method for weighted networks is called the *weight rank clique filtration* (WRCF), which filters the network based on weights. The procedure works as follows:

1. Set the set of all nodes, without any edge, as filtration step 0.
2. Rank all edge weights in decreasing order $\{w_1, \dots, w_n\}$.
3. At filtration step t , keep only the edges whose weights are less than w_t , thus creating an unweighted graph.
4. Define the maximal cliques of the resulting graph to be simplices.

At each step of the filtration, we construct a simplicial complex based on cliques: this is called a *clique complex*. It is necessarily valid since a subset of a clique is necessarily a clique itself, and the same is true for the intersection of two cliques.

This leads to a first possibility for applying persistent homology to temporal networks. It is possible to segment the lifetime of the network into sliding windows, creating a static graph on each window by retaining only the edges available during the time interval. We can then apply WRCF on each static graph in the sequence, obtaining a filtered complex for each window, to which we can then apply persistent homology.

This method can quickly become very computationally expensive, as finding all maximal cliques (using the Bron-Kerbosch algorithm for example) is a complicated problem in itself. In practice, we often restrict the search to cliques of dimension lower than a certain bound d_M . With this restriction, the new simplicial complex is homologically equivalent to the original: they have the same homology groups up to H_{d_M-1} .

This method is sensitive to the choice of sliding windows on the time scale. The width and the overlap of the windows can completely change the networks created and their topological features. Too small a window, and the network becomes too small to have any significant topological properties, too large, and we lose important information in the evolution of the network over time.

3.4 ZIGZAG PERSISTENCE