

Promises and pitfalls of topological data analysis for brain connectivity analysis[☆]



Luigi Caputi^{a,1}, Anna Pidnebesna^{a,b,c}, Jaroslav Hlinka^{a,b,*}

^a Institute of Computer Science of the Czech Academy of Sciences, Pod Vodárenskou věží 271/2, Prague 182 07, Czech Republic

^b National Institute of Mental Health, Topolová 748, Klecany 250 67, Czech Republic

^c Faculty of Electrical Engineering, Czech Technical University, Technická 1902/2, Prague 166 27, Czech Republic

ARTICLE INFO

Keywords:

Persistent homology
Connectivity
fMRI
Electrophysiology
Epilepsy
Schizophrenia

ABSTRACT

Developing sensitive and reliable methods to distinguish normal and abnormal brain states is a key neuroscientific challenge. Topological Data Analysis, despite its relative novelty, already generated many promising applications, including in neuroscience. We conjecture its prominent tool of persistent homology may benefit from going beyond analysing structural and functional connectivity to effective connectivity graphs capturing the direct causal interactions or information flows. Therefore, we assess the potential of persistent homology to directed brain network analysis by testing its discriminatory power in two distinctive examples of disease-related brain connectivity alterations: epilepsy and schizophrenia. We estimate connectivity from functional magnetic resonance imaging and electrophysiology data, employ Persistent Homology and quantify its ability to distinguish healthy from diseased brain states by applying a support vector machine to features quantifying persistent homology structure.

We show how this novel approach compares to classification using standard undirected approaches and original connectivity matrices. In the schizophrenia classification, topological data analysis generally performs close to random, while classifications from raw connectivity perform substantially better; potentially due to topographical, rather than topological, specificity of the differences. In the easier task of seizure discrimination from scalp electroencephalography data, classification based on persistent homology features generally reached comparable performance to using raw connectivity, albeit with typically smaller accuracies obtained for the directed (effective) connectivity compared to the undirected (functional) connectivity. Specific applications for topological data analysis may open when direct comparison of connectivity matrices is unsuitable - such as for intracranial electrophysiology with individual number and location of measurements. While standard homology performed overall better than directed homology, this could be due to notorious technical problems of accurate effective connectivity estimation.

1. Introduction

The general idea that mathematical approaches from geometry and algebraic topology can provide powerful tools and methods for the understanding of complex data is fairly recent (see, e.g., (Carlsson, 2009; Edelsbrunner et al., 2002; Frosini, 1990; Zomorodian and Carlsson, 2005)) and has led to the birth and advancement of a completely new research branch, *Topological Data Analysis* (TDA). This is a fast growing and promising field at the intersection of Data Analysis, Algebraic Topology, Computational Geometry, Machine Learning and Statis-

tics. Its main advantage is the possibility to analyse data by looking at its shape and underlying geometric structures arising from high-dimensional relations; this goal is achieved by providing a strong mathematical framework and manageable quantitative tools such as homology theory (Hatcher, 2000; Munkres, 1984). The structure of the data is then qualitatively and quantitatively assessed in the form of topological features (e.g. connected components, voids, tunnels or loops). This approach has been already widely applied and finds nowadays applications in many science domains dealing with various complex systems (Chazal and Michel, 2017; Wasserman, 2018).

[☆] The study was supported by the Czech Science Foundation projects No. GA19-11753S, GA21-14727K, GA21-17211S, GA21-32608S, by the Czech Health Research Council projects No. NV17-28427A and NU21-08-00432, and by Ministry of Health of the Czech Republic – DRO 2021 („National Institute of Mental Health – NIMH, IN: 00023752“).

^{*} Corresponding author at: Institute of Computer Science of the Czech Academy of Sciences, Pod Vodárenskou věží 271/2, 182 07 Prague, Czech Republic.
E-mail addresses: luigi.caputi@abdn.ac.uk (L. Caputi), pidnebesna@cs.cas.cz (A. Pidnebesna), hlinka@cs.cas.cz (J. Hlinka).

¹ Present address: Institute of Mathematics, University of Aberdeen, AB24 3UE, UK.

<https://doi.org/10.1016/j.neuroimage.2021.118245>.

Received 21 January 2021; Received in revised form 30 May 2021; Accepted 5 June 2021

Available online 7 June 2021.

1053-8119/© 2021 The Authors. Published by Elsevier Inc. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>)

The human brain and its dynamical activity represents one of the archetypal examples of complex systems. Indeed, due to both the intricate nature of the structure of connections of its constituent units, as well as the highly nonlinear and potentially chaotic dynamics within each of these units, it is rightly considered a suitable target example system for application of such complex system methodology. In particular, the study of the complex patterns of brain network topology has become a flourishing area of research since the 2000/s (Bullmore and Sporns, 2009).

There are three main approaches to assess the interactions between brain regions: structural, functional, and effective connectivity (Friston, 1994; Friston, 2011; Park and Friston, 2013). Structural connectivity is the pattern of physical connections between the considered brain regions, while functional and effective are based on neuronal activity signals. Functional connectivity is defined as the stochastic dependence between remote neurophysiological events and provides thus a non-directional (symmetric) measure of brain interactions. On the contrary, the effective connectivity is defined as the direct influences that one region exerts over another.

The most frequently used measure of functional connectivity is the (linear) Pearson correlation, which provides an accurate functional connectivity estimate for functional magnetic resonance imaging data under common processing settings (Hlinka et al., 2011, Hartman et al., 2011). One of the applications in clinical neuroscience is the characterization of functional brain alterations in various diseases. The classification power of FC has been reported for schizophrenia (Lynall et al., 2010), multiple sclerosis (Zurita et al., 2018), epilepsy (Liao et al., 2010), and many other brain diseases and states. On the other hand, effective connectivity, while in principle more informative with respect to the underlying dynamics since it specifies the direct influences among the brain regions, is more scarcely applied in practice, presumably since it is also much harder to estimate and the methods thus typically can handle only smaller networks (Smith, 2012). While there are many competing methods for estimation of effective connectivity offering different trade-off between advantages and drawbacks (Valdes-Sosa et al., 2011), there is a general agreement that if estimated well, they would provide valuable extra information on top of functional connectivity (Smith, 2012).

While a connectivity matrix characterizes only the strengths of all pair-wise interactions, more complex characterizations of the whole network as well as specific nodes thereof can be provided. This can be done by application of graph-theoretical tools to the graph constructed by using the connectivity matrix as an adjacency matrix (Telesford et al., 2011). However, the graph-theoretical analysis of brain connectivity suffers from multiple technical difficulties. One of them is that the classical approach requires applying a threshold to select which connections are relevant, i.e., strong enough to constitute an edge in the graph. There are multiple heuristics to deal with the threshold selection or even generalization of many graph metrics to weighted graphs, however no single optimal solution is available. Among other attempts Gallos et al. (2012), a very attractive way to side-step the problem and consider all possible thresholds at once, without losing any information, is thus offered by already established methods of TDA, such as Persistent Homology.

Persistent Homology (PH) is one of the main tools adopted in TDA and is a multi-scale adaptation of the classical homology theories studied in Algebraic Topology. It allows a computation of the *persistent* topological features of a space at all different resolutions (here: densities of the thresholded connectivity graph), while revealing the most essential ones. In the PH-approach to complex networks, instead of focusing on a single thresholded graph, one analyses a family of simplicial complexes indexed by a (filtration) parameter, by considering all the possible thresholds at once.

Based on a robust mathematical theory and being stable with respect to small noise perturbations (Cohen-Steiner et al., 2007), PH has provided a novel qualitative and quantitative tool to study complex data (in the form of either point clouds, time series or connectivity networks). Its use has already been reported in many different fields: in neuroscience

and neurological disorders (Lee et al., 2011b), in endoscopy analysis (Dunaeva et al., 2016), angiography (Bendich et al., 2014), pulmonary diseases (Brodzki et al., 2018), finance (Gidea, 2017), fingerprint classification (Giansiracusa et al., 2019), image classification (Dey et al., 2017), to name few applications. The pipeline involving Persistent Homology in network analysis usually starts by considering (undirected, weighted) graphs as input. However, in many real world examples as in the study of effective connectivity, the input graphs are also equipped with a richer structure on the edges, namely the direction from the source node to the target. As directed graphs theoretically hold more information than their undirected counterparts (see e.g., Section 3.4), an extension of the PH pipeline to asymmetric networks might provide a more refined family of invariants. Some novel approaches have already been proposed in this direction by Masulli and Villa (2016), Reimann et al. (2017), Chowdhury and Mémoli (2018), Turner (2019), Aktas et al. (2019), Lütgehetmann et al. (2020), although not yet widely adopted; notably the early applications include the analysis of brain structural connectivity networks (Reimann et al., 2017).

The developments in the fields of brain connectivity analysis and TDA call for their combined utilization. Many attempts have been done in this direction, for reviews see (Lee, 2019; Phinyomark et al., 2017; Solo et al., 2018). In this work, with the aim of investigating the promise of applications of topological analysis to effective connectivity networks, we provide a pioneering example of such endeavour for the task of classifying altered brain states in several settings, while we assess the improvement provided by these methodologies in contrast to dealing with the connectivity directly, or applying the TDA approach to functional connectivity data only.

To this end, we construct functional and effective connectivity matrices (in short, FC and EC) derived from brain activity time series; these are obtained by two most common functional neuroimaging modalities – (intracranial) electroencephalography ((i)EEG) and functional magnetic resonance imaging (fMRI) – using two prominent examples of disconnectivity-related brain diseases: epilepsy and schizophrenia. For epilepsy, for a group of patients, we distinguish between ictal and pre-ictal connectivities. For schizophrenia, we classify patients from the matching neurotypical controls. Following previous approaches Lee et al. (2011b, 2012); Merelli et al. (2016); Stolz et al. (2018), we calculate the persistent homology of the (directed) filtrations associated to the brain functional (or effective) connectivity networks, and compute the associated PH-features (further referred to as *directed persistent homology*). In order to quantify their informative power, we apply a state-of-art classifier (a support vector machine, SVM), with the task to distinguish between healthy and diseased brain states. We use a range of most common quantitative PH features, namely *Persistence Landscapes* (PL) Bubenik and Dłotko (2017), *Persistence Images* (PI) Adams et al. (2017), *Persistent Entropy* (PE) Rucco et al. (2014) and *Carlsson Coordinates* (CC) Adcock et al. (2013). We compare the results with those obtained, under the same tasks, by a naive (SVM) classification applied to the FC and EC matrices directly.

The directed homology pipeline followed in this work was introduced in Reimann et al. (2017) with applications in the study of the structural and functional connectivity in rat brain slices. This was also the first, and the only one known to the authors, application of this approach in neuroscience.

Note that while the application of directed persistent homology analysis to human brain effective connectivity is novel, it builds upon the tools and ideas developed and applied recently in the field. For example, the undirected variants of the TDA features have been shown to be promising in applications to Alzheimer (Kuang et al., 2019b; Pachauri et al., 2011), ASD (Wong et al., 2016), depression (Khalid et al., 2014), Parkinson (Garg et al., 2016), or in the study of childhood stress (Chung et al., 2014). The first work, applying PH-methods to brain networks constructed by using Pearson correlation is due to Lee (see Lee et al., 2011b; 2012) in the case of patients affected by ADHD. Among previous approaches in epilepsy, a persistent homology analysis has

been already reported in (Merelli et al., 2016), where (weighted) PE is employed in order to detect the transition between the pre-ictal and ictal states, or in (Piangerelli et al., 2018; Wang et al., 2015), where a classification analysis in epilepsy, based on PLs and PE, has used different filtrations (sublevel-based or lower star filtrations). In (Choi et al., 2014), a study of the epileptic abnormality is given in rat models. Within schizophrenia applications, the only analysis applying methods from TDA is, to the authors knowledge, the work (Stolz et al., 2018). To provide more insight into the promises and pitfalls of the application of directed persistent homology analysis of brain (effective) connectivity, the current study goes beyond a single method application, and provides comparisons with the application of (standard) persistent homology of undirected networks, as well as with the application of machine learning directly to the input connectivity matrices. This, together with tackling three types of neuroimaging data across two brain different dysconnectivity conditions, namely through an easier classification task of ictal versus pre-ictal brain states, and a more compelling one, between schizophrenia patients and healthy controls, should provide a richer picture concerning the challenges of TDA application to brain connectivity data analysis.

Our key observation is that standard TDA applied to effective connectivity is in principle able to discriminate diseased from healthy brain states, comparably to (previously used combinations of) TDA applied to functional connectivity. However, in many practical situations such as we demonstrate in the results section for the application to schizophrenia classification, the use of raw functional connectivity features is advantageous. This may be due to topographical, rather than topological, specificity of the particular changes related to schizophrenia – these topographically (regionally/spatially) specific changes are naturally picked up by classifiers using the raw functional connectivity matrices; while the topological methods effectively ignore the spatial information, as they are invariant with respect to permutation of nodes – a fact that may be commonly overlooked when considering the applications. However, in yet another scenarios, the direct connectivity comparison is not straightforwardly applicable, such as for inter-subject analysis of intracranial electrophysiology recordings in epilepsy, not only due to potentially different spatial location of the brain alterations, but also due to the technical challenge of inter-individual differences in the number and location of the measurement sites.

With the idea that the topological pipeline may suffer from heterogeneous subject biases, a probable cause of lower general performances, we have also introduced an alternative methodological pipeline reducing the underlying heterogeneity. We observe that this variation of the pipeline shows a remarkable increment of the performances in the TDA-approach, when applied to the EEG dataset.

The paper is structured as follows. In Section 2 we recall the definition of directed Persistent Homology. In Section 3, we describe our pipeline and the common featureizations adopted in the classification tasks; we report our results in Section 4. In Section 5, we discuss our contributions. In the Appendix A, we review some mathematical concepts in more detail, including simplicial complexes, filtrations and (persistent) homology groups. Appendix B and Appendix C provide additional details concerning the data selection and processing. We conclude with some additional results (Appendix D) and a statistical comparison of PH versus DPH (Appendix F).

2. Persistent homology of directed networks

In this section, following Reimann et al. (2017) and Lütgehetmann et al. (2020), we review the classical notions of simplicial complexes and persistent homology groups adapted to the *directed* framework, keeping always in mind the directed networks as the application example. The main objects in this context are directed graphs, and as in the context of undirected graphs, starting from a directed network we will construct suitable simplicial complexes, filtrations and persistent homology groups. We assume that the reader

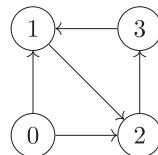
is familiar with these classical general concepts in the undirected case (reviewed for convenience in Appendix A); we refer to Zomorodian and Carlsson (2005) and to Otter et al. (2017) for an introduction of Persistent Homology with an eye towards applications.

A graph $G = (V, E)$, where V denotes the set of its vertices and E the set of its edges, is a *directed graph* if every edge $e = (v_0, v_1)$ is given by an *ordered* pair of *distinct* vertices v_0 (the source of the edge) and v_1 (the target). For a given pair of vertices v_0 and v_1 , both edges (v_0, v_1) and (v_1, v_0) are allowed, but loops, i.e., edges of type (v, v) , are not. A higher dimensional generalization of directed graphs is given by the notion of directed simplicial complexes (see Def. Appendix A.1 for the classical definition of simplicial complexes):

Definition 2.1. An *ordered simplicial complex* $\Sigma = \{\sigma_\alpha\}_\alpha$ on a vertex set V is a non-empty family of finite ordered subsets $\sigma_\alpha \subseteq V$ with the property that, if σ belongs to Σ then every ordered subset τ of σ (ordered with the natural order induced by σ) belongs to Σ .

In analogy with the notation used for simplicial complexes, the elements of an ordered simplicial complex of cardinality $d + 1$ are called *d-simplices*. We say that an ordered simplicial complex Σ has dimension d if all its simplices have cardinality at most $d + 1$. For example, every directed graph (V, E) is a 1-dimensional ordered simplicial complex on the set of vertices V and, if V is the set $\{0, \dots, n\}$, then the collection of all the ordered subsets of V is an n -dimensional ordered simplicial complex. When $n = 1$ and $V = \{0, 1\}$, the ordered simplicial complex given by the ordered subsets of V is nothing but the directed graph with vertices 0 and 1 and ordered edges $(0,1)$ and $(1,0)$.

In the context of undirected graphs, there are various ways to construct associated simplicial complexes; one of the most established ones is given by the clique complex, where one builds the simplicial complex associated to a graph G by using the combinatorial information of its cliques (see Ex. Appendix A.3). Following the analogy with this classical construction, from a given directed graph G we will construct an ordered simplicial complex by using the combinatorial information given by the ordered cliques of G , as we now explain. Recall that a k -clique in an unordered graph G is a complete subgraph on k vertices. Likewise, if $G = (V, E)$ is a directed graph, then an *ordered k-clique* of G is a k -tuple (v_1, \dots, v_k) of vertices of G with the property that the pair (v_i, v_j) is an ordered edge of G for every $i < j$. We observe that in an ordered clique there is a source vertex, v_1 , having only out-edges, and a target vertex, v_k , having only in-edges. For example, in the following graph on $V = \{0, 1, 2, 3\}$:

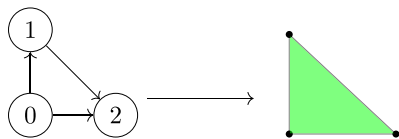


we have only one ordered 3-clique (in addition to the 1- and 2-cliques), i.e., the clique corresponding to the complete subgraph $(0,1,2)$ on the vertices 0, 1 and 2. Observe that the subgraph on the vertices 1, 2 and 3, although being a 3-clique for the underlying undirected graph, is not an *ordered* 3-clique as there is no common source or target (it is, in fact, a cycle). For a given directed graph, we can construct an ordered simplicial complex by considering the collection of all the directed cliques (see the analogous construction of Ex. Appendix A.3):

Definition 2.2. The *directed flag complex* $dFl(G)$ associated to a directed graph G is the ordered simplicial complex whose k -simplices are all the ordered $(k + 1)$ -cliques of G .

Observe that the vertices and edges of a directed graph G are exactly the 1- and 2-cliques of G , and that every directed graph is contained in its associated directed flag complex. For example, if G is the directed graph on $V = \{0, 1, 2\}$, with edges $(0, 1)$, $(1, 2)$ and $(0, 2)$, then there is only one directed 3-clique, $(0,1,2)$, and the associated directed flag complex

$dFl(G)$ consists of the same vertices and edges of G , together with the 2-simplex $(0,1,2)$ corresponding to the unique directed clique:



In order to study persistent homology groups, once we have a suitable notion of simplicial complexes, we need filtrations of simplicial complexes. A filtration of ordered simplicial complexes consists of a nested sequence of increasing ordered simplicial complexes (see also Def. Appendix A.4). Various filtrations are generally possible. In Ex. Appendix A.5 we review the clique filtration and in Ex. Appendix A.6 the Rips filtration, an instance of the clique filtration (see Remark Appendix A.7) for undirected graphs. In the context of directed graphs, we construct filtrations using the directed flag complexes, as the natural extension of the clique filtration to the directed context. To be more precise, if $\emptyset \subseteq G_0 \subseteq \dots \subseteq G_n = G$ is a family of directed subgraphs of a directed graph G , then the sequence

$$\emptyset \subseteq dFl(G_0) \subseteq \dots \subseteq dFl(G_n) = dFl(G) \quad (2.1)$$

of directed flag complexes $dFl(G_i)$ is a filtration of $dFl(G)$.

In our applications, we first consider a *weighted* directed graph G , i.e., a directed graph for which every edge is given a numerical weight. For each weight w_i we define the directed graph G_i as the (directed) subgraph of G with same vertices of G and only those edges corresponding to the weights w of G with $w \leq w_i$ (note that at this level the weights correspond already to a distance, rather than similarity matrix, see also Ex. Appendix A.5). The weights w_i of G , sorted in the increasing order, induce a sequence of directed subgraphs of G , hence a filtration of directed flag complexes $dFl(G_i)$. Generally speaking, in order to extend it to a filtration of directed flag complexes one has to choose a weighting function also on the directed simplices (hence on the ordered cliques) extending the weights given on the edges. We use the intuitive maximum weighting function (the default in the software *Flagser* Lütgehetmann et al. (2020)), associating to an ordered simplex the maximal weight across all its edges. See Lütgehetmann et al. (2020) or the documentation at <https://github.com/luetge/flagser> for details on other possible weighting functions.

Once a filtration of (ordered) simplicial complexes is given, one can apply the classical homology methods for obtaining Persistent Homology groups, barcodes and Persistence Diagrams (see Appendix A.2 for a review). We will briefly refer to the Persistent Homology of the filtration of directed flag complexes associated to a weighted directed graph as its *directed Persistent Homology* (shortly denoted by DPH).

3. Methods

In the analysis of brain connectivity networks, we use Persistent Homology features as input for classification by an SVM classifier. The steps of the pipeline include (see also Fig. 1):

- Data:** we use scalp and intracranial EEG (of epileptic patients) and fMRI (of schizophrenia patients and healthy controls) data in the form of multi-variate time series;
- Graphs:** for each dataset we construct (functional and effective) connectivity graphs by calculating Pearson correlation and Granger Causality respectively; see Section 3.2.
- Topology:** from the graphs we construct appropriate filtrations of simplicial complexes (see Section A.1 and Section 2). In particular, we have constructed the Vietoris-Rips filtration (see Ex. Appendix A.6) associated to the FC graph and the filtration of directed flag complexes associated to the directed EC graph (see Eq. (2.1)).
- Homology:** we compute the Persistent Homology of the given filtrations (PH for undirected and DPH for directed graph), getting

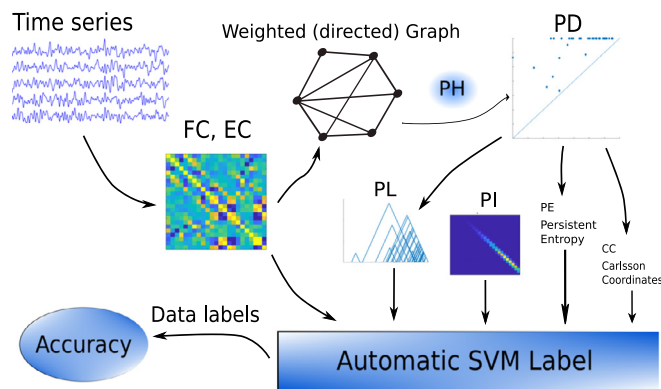


Fig. 1. Pipeline of the main analysis. From the time series, we estimate the connectivity matrix, based on which build a weighted graph. After estimating Persistent Homology (PH) of the graph, we get Persistent Diagram (PD), that is represented by different features (Persistent Landscape PL, Persistent Image PI, Persistent Entropy PE, Carlsson Coordinates CC) for classification. Finally, classification based on PH features is compared with classification using matrices itself.

their persistence Betti numbers (corresponding to the number of connected components for dimension 0 and to the number of loops for dimension 1, the two most common settings in the literature); see also Appendix A.2.

- Machine Learning:** from persistence Betti numbers, we get the features (vectors of real numbers) and apply a suitable classifier (SVM) to both the FC or EC graphs, and to the features (in particular, the Persistence Landscapes, Persistent Images, Carlsson Coordinates and Persistent Entropy (Section 3.5)).

As a variation of this pipeline, we add an additional cleaning step in the case of EEG and iEEG datasets as described in Section 3.6.

3.1. Dataset

We use two prominent types of neuroimaging data in combination with two examples of pathological brain dynamics:

3.1.1. fMRI data (Schizophrenia)

We analyse the dataset that consists of fMRI recordings of 90 patients with a schizophrenia diagnosis and 90 healthy controls. Functional MRI data were collected with a 3T MR scanner (Siemens; Magnetom Trio) at Institute of Clinical and Experimental Medicine in Prague. Participants were informed about the experimental procedures and provided written informed consent. The study design was approved by the local Ethics Committee of the Institute of Clinical and Experimental Medicine and the Psychiatric Center Prague. $T2^*$ -weighted images with BOLD contrast were acquired with voxel size $3 \times 3 \times 3$ mm³ and TR/TE parameters of 2000/30ms. As anatomical reference the T1-weighted high-resolution image was used (TR/TE/TI = 2300/4.6/900 ms, voxel size = $1 \times 1 \times 1$ mm³). Initial data preprocessing was done using FSL routines (FMRIB Software Library v5.0, Analysis Group, FMRIB, Oxford, UK) and CONN toolbox (McGovern Institute for Brain Research, MIT, USA). Functional realignment and unwarping, slice-timing correction, structural segmentation into white matter and cerebrospinal fluid and structural normalization to the MNI space were done using CONN's default preprocessing pipeline (defaultMNI), as well as functional normalization to the MNI space, outlier detection, and smoothing with 8mm kernel size. The estimated amount of motion in the subjects was generally satisfactory, with the average framewise displacement between subsequent volumes being 0.13 ± 0.08 mm (mean \pm std across subjects). Six head-motion parameters with their first order temporal derivatives and five principal components of white-matter and cerebrospinal fluid time-series were regressed out

to suppress measurement artifacts. Furthermore, to remove possible signal drift, time-series were linearly detrended and filtered by band-pass filter [0.009-0.08Hz]. See Kopal et al. (2020) and Oliver et al. (2019) for detailed preprocessing description. To extract the time series for further analysis, the brains spatial domain was divided into 90 non-overlapping regions of interest (ROIs) according to the AAL atlas; from each ROI we extract one BOLD time series by averaging the time series of all voxels in the ROI.

3.1.2. EEG data (Epilepsy)

We analyse both scalp and intracranial electroencephalography EEG data:

Scalp EEG data. We use EEG recordings from pediatric subjects with intractable seizures, collected at the Children's Hospital Boston. The data are available at <https://physionet.org/content/chbmit/1.0.0/>, from the publicly open dataset *Physionet*. We analyse EEGs belonging to 18 subjects, with 23 scalp electrodes. In particular, we included patients with typical epileptic event of length between 20 and 120 s. The set of selected patients is provided in [Appendix B](#). After isolating the epileptic event according to labels in the patient's summary file, we also consider non-epileptic records, chosen randomly from the database (we extract segments of the same length as the seizure segments). In total, we use 102 ictal and 102 interictal data segments, the same amount of ictal and preictal parts for every particular subject. Sampling frequency is 256Hz, no special preprocessing is applied except for a band-pass filter [1-70Hz], notch filter at 50Hz, and removal of global signal for every data segment separately.

Several seizure recordings are available for each subject. Therefore, to avoid dependence between samples and any potential biases due to different numbers of seizure recordings available, for each subject we average all the connectivity matrices for ictal (interictal) periods. Thus, the dataset for the analysis consists of two connectivity matrices per subject (one ictal and one inter-ictal). For the additional analysis using a modified pipeline ([Section 3.6](#)) we used a slightly modified data selection, which is described in [Appendix B](#).

Intracranial EEG data. We use the anonymized dataset of 16 patients of the epilepsy surgery program of the Inselspital Bern, freely available at <http://ieeg-swez.ethz.ch/>, with no further preprocessing implemented. The intracranial electroencephalography (iEEG) recordings belong to 16 patients with pharmacoresistant epilepsies collected at the Sleep-Wake-Epilepsy-Center (SWEC) of the University Department of Neurology at the *Inselspital Bern* and the *Integrated Systems Laboratory* of the ETH Zurich. Each recording consists of 3 min of preictal segments, the ictal segment (ranging from 10 s to 1002 s), and 3 min of postictal time. In this work we only consider the recordings corresponding to the first 30 s of the first seizure as the 'ictal' data, and the 30 s segment starting 1 min before the first seizure start as our 'preictal' data.

For every patient we use the first minute of interictal data for estimation of baseline connectivity, which is later subtracted from both ictal and preictal connectivities.

3.2. Connectivity networks

In order to estimate the brain networks we use both a symmetric (Functional Connectivity, [Section 3.2.1](#)) and an asymmetric approach (Effective Connectivity, [Section 3.2.2](#)):

3.2.1. Functional connectivity

To obtain the functional connectivity network, we cross-correlate the time-series described in [Section 3.1.2](#) and [Section 3.1.1](#). Let N be the number of brain regions in the AAL atlas or number of electrodes in scalp EEG, then the $N \times N$ matrix $A = (A_{i,j})$ is defined using Pearson correlation between the time series corresponding to nodes i and j . In line with previous literature (Lee et al., 2011b, [Section 3](#)), (Merelli et al., 2016, [Section 3](#)), (Stolz et al., 2018, [Section 2.2](#)), all entries of A that are not significant at $\alpha = 0.05$ are set to 0, along with all negative values

of functional connectivity. Below, we refer to this method as *masked Functional Connectivity (masked FC)*.

3.2.2. Effective connectivity

The effective connectivity is quantified by the pairwise Granger causality (Barnett and Seth, 2013); furthermore, edges corresponding to links that do not show statistically significant causal effect (by means of the FACDA algorithm (Kořenek and Hlinka, 2020)) have their weight set to zero. In the following we refer to this method as *masked Effective Connectivity*, or more concisely, as *masked EC*. We use pairwise Granger causality instead of multivariate Granger Causality (which controls for all the other variables) to reach more robust estimates. Indeed, in the multivariate approach a greater number of parameters has to be estimated, potentially leading to less accurate estimates (Hlinka et al., 2013).

3.3. Persistent homology computation

We analyse persistent homology groups of both the directed and undirected networks. We refer to [Appendix A](#) for a review of PH. When not otherwise specified, the persistent homology computations shown in the article refer to calculations implemented in dimension 1, that is detection 1-dimensional holes; we refer to [Appendix D](#) for the corresponding results using the dimension 0 features, corresponding to the detection of connected components; homology groups are computed over the commonly used base field \mathbb{F}_2 (see Def. [Appendix A.8](#)).

3.3.1. PH of undirected networks

The undirected network constructed in [Section 3.2](#) is a weighted graph. As explained in Ex. [Appendix A.5](#), by sorting the weights in an increasing or decreasing order, we get a filtration of simplicial complexes, the WRCF. Equivalently, one can construct the associated Vietoris-Rips filtrations (see Ex. [Appendix A.6](#)). As the implementation of the WRCF can be computationally demanding, in this work we follow the second approach: we construct the quasi-metric space (X, d) associated to the matrix A (the matrix A being constructed as described in the previous subsection) with underlying space X given by the vertices of the weighted graph and quasi-distance $d(i, j) := 1 - A(i, j)$. We then construct the Vietoris-Rips filtration associated to this metric space and compute the associated PH and persistence diagrams. The computation of the associated Persistent Homology has then been implemented with the python library *Ripsper* (Tralie et al., 2018).

3.3.2. PH of directed networks

The directed network presented in [Section 3.2.2](#) provides a weighted directed graph, hence a directed flag complex as explained in [Section 2](#). For consistency with the undirected case, we first rescale the matrix A via a linear map $A \mapsto \tilde{A}$ to span the interval $[0, 1]$ and then apply the transformation $x \mapsto 1 - x$ to use the matrix $1 - \tilde{A}$, getting a positive matrix where the lowest values represent the strongest connections. We consider the sequence of directed graphs sorted with the increasing order of the weights (of $1 - \tilde{A}$). We extend the weights to the simplices of the induced flag complexes by giving to a simplex σ the maximum value across the weights of its edges. We then compute the associated persistent homology groups.

3.4. Why DPH? A motivational example with VAR processes

In this subsection we elucidate why an extension of the persistent homology techniques to directed networks is theoretically more fruitful than the classical PH approach on undirected networks. We assume that the reader is familiar with the concept of VAR(p) (vector autoregressive) processes (see Lutkepohl (2007) for a general overview).

Let X be a VAR(1) process of the form

$$X_t = AX_{t-1} + e_t$$

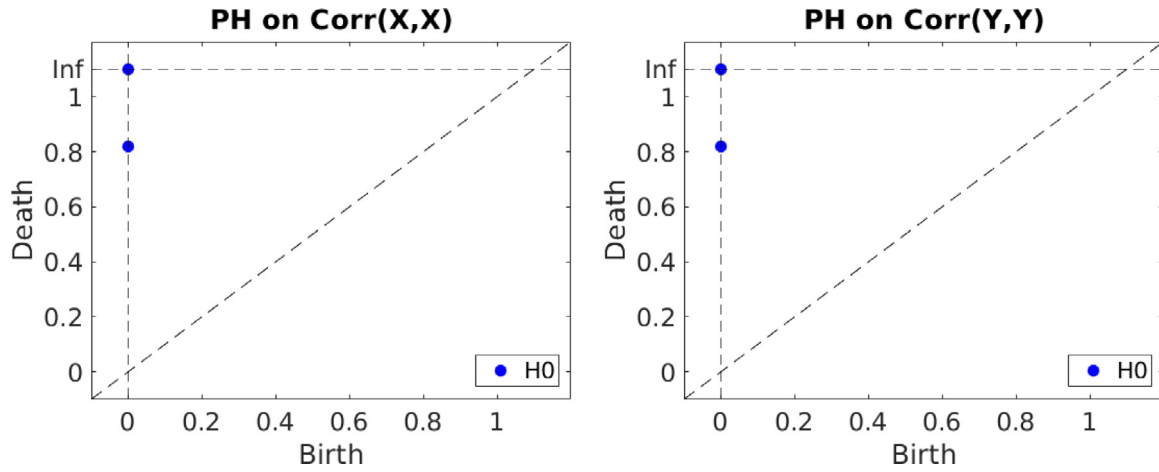


Fig. 2. Left: PH of correlation matrix $Corr(X, X)$. Right: PH of correlation matrix $Corr(Y, Y)$.

where X_t is an n -dimensional random vector, A is a $n \times n$ matrix of (real) coefficients and e_t is a white noise with 0 mean and covariance Σ ; we say that the process X is *generated* by the matrix A .

Without loss of generality, we assume that X is *stable* (meaning that the matrix A has eigenvalues smaller in module than 1), has 0 mean and the variance Σ is the identity matrix $\Sigma = I$. Then, the matrix sequence $\{A^i\}_{i \in \mathbb{N}}$ is absolutely summable and the covariance matrix of X (that is a stable process, hence stationary) can be written as

$$\text{Var}(X) = (I - AA^t)^{-1} \tag{3.1}$$

where A^t denotes the transpose of the matrix A .

The correlation matrix of the process X depends only on its covariance matrix $\text{Var}(X) = \text{Cov}(X, X)$, whereas the associated effective connectivity depends upon the generating matrix A (for simplicity, we use the matrix A directly here as the EC).

In order to document the extra richness of the directed version of PH, we consider two different VAR(1) processes X and Y generated by two matrices A and B with different topological structure but the same covariance matrix $\text{Cov}(X, X) = \text{Cov}(Y, Y)$. These two processes share the same correlation matrix (used in the definition of FC), hence the associated (standard/undirected) PH will not distinguish them. However, as the associated EC matrices differ, the *directed* PH might be able to distinguish the aforementioned processes. We now document this using a simple low dimensional example.

The problem reduces to finding topologically different matrices A, B , such that $AA^t = BB^t$ (this can be achieved by choosing an suitable orthogonal matrix Q , and by setting $B := AQ$. As a specific example, we set:

$$A := \frac{1}{4} \begin{pmatrix} 0 & 1 & 2 \\ 1 & 0 & 1 \\ 2 & 1 & 0 \end{pmatrix}, Q_\varphi := \begin{pmatrix} \cos(\varphi) & \sin(\varphi) & 0 \\ -\sin(\varphi) & \cos(\varphi) & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

Let $\varphi = \pi/3$ and $B = AQ_\varphi$. The eigenvalues of A and B are smaller in module than 1 and the generated VAR(1) processes are stable, moreover, clearly $AA^t = BB^t$, hence, they have the same covariance matrix by Eq. (3.1), and as a consequence, the correlation matrices of X and Y have same persistent homology groups (see Fig. 2). On the other hand, when applied to the generating matrices A and B , DPH is able to distinguish the two processes (see Fig. 3), proving the advantage of DPH in distinguishing a larger class of processes.

We point out here that we cannot apply PH directly on B , as B is an asymmetric matrix. However, PH could indeed be in principle applied to the matrix A that is symmetric; however, it provides a different homology structure than the application of DPH. In particular, the persistent homology groups associated to A are concentrated in dimension 0, whereas DPH(A) has non-trivial homology classes in dimension 1 and 2 as well as in dimension 0 (see Fig. 4).

3.5. Featurizations of persistent homology

Persistent homology groups can be visualized by using Persistence Diagrams as birth-death planar representations (Section A.2). In order to apply a SVM classifier, we derive four most commonly used vectorizations, or featurizations, from the persistence diagrams: Persistence Landscapes, Persistence Images, Carlsson Coordinates and Persistent Entropy. In this way we get input real-valued vectors amenable to the classifier. We proceed by illustrating the definitions of these topological features.

3.5.1. Persistence landscapes (PL)

Persistence landscapes have been introduced in order to perform PH-based statistical analysis by Bubenik and Dłotko (2017). For a birth-death interval (b, d) , where d is a finite value, the function $f_{(b,d)} : \mathbb{R} \rightarrow \mathbb{R} \quad t \mapsto \max(\min(t - b, d - t), 0)$ represents a triangle-shaped function over (b, d) . Let $\{(b_i, d_i)\}_{i \in I}$ be a (finite) collection of birth-death intervals. Let $k - \max(\{a_i\}_i)$ denote the k th maximum value of a finite family of real numbers $\{a_i\}_i$.

Definition 3.1. The *persistence landscape* of the set of birth-death coordinates $\{(b_i, d_i)\}_{i \in I}$ is the family of functions $\lambda = (\lambda_k : \mathbb{R} \rightarrow \mathbb{R})_{k \in \mathbb{N}}$, where $\lambda_k(t) := k - \max_{i \in I} (f_{(b_i, d_i)}(t))$.

Observe that the functions $\lambda_k(t)$ are piece-wise linear and that they are identically 0 when k is bigger than the cardinality $|I|$ of the birth-death pairs. See Fig. 5 for an example of persistence landscape associated to a persistence diagram.

We follow the same method as Yesilli et al. (2019) to perform ML analysis on the set of persistence landscapes. Let $D = (D_p)_{p \in P}$ be a set of persistence diagrams (e.g. each diagram D_p is the PD associated to a single brain network, at a fixed dimension). For every such persistence diagram D_p , let $\lambda^{(p)} = \{\lambda_i^{(p)}\}_{i \in \mathbb{N}}$ be the associated PL. Let $K \subseteq \mathbb{N}$ be a fixed subset of indices of \mathbb{N} ; for every $k \in K$ consider the collection of functions $\{\lambda_k^{(p)}\}_p$ varying with the index p (and fixed k).

By definition, for every k and p , the map $\lambda_k^{(p)} : \mathbb{R} \rightarrow \mathbb{R}$ is a piecewise linear function. Then, it has a set of critical points $\{x_{k,i}^{(p)}\}_i$, i.e., those points at which the function $\lambda_k^{(p)}$ changes its direction. Let $z^{(k)}$ be the vector of \mathbb{R}^n whose entries consist of all the critical points $\{x_{k,i}^{(p)}\}_{i,p}$ associated to the functions $\lambda_k^{(p)}$, at k fixed; the entries of $z^{(k)}$ are sorted in increasing order, without repetitions.

For every chosen $k \in K$, we have constructed a vector $z^{(k)}$ depending only on the critical points of the k th functions of the PLs. The feature vector, input of the classifier, associated to the diagram D_p in degree k , is then the evaluation vector $\{\lambda_k^{(p)}(z_j^{(k)})\}_{0 \leq j \leq |z^{(k)}|}$. See the work of Chazal and Michel (2017) for further details. In our work, we choose

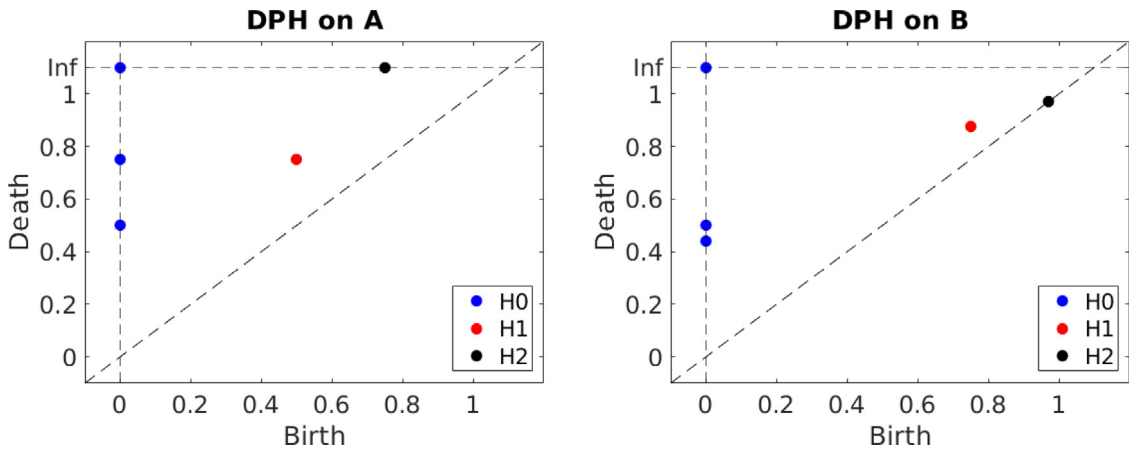


Fig. 3. Left: DPH of the matrix A. Right: DPH of the matrix B.

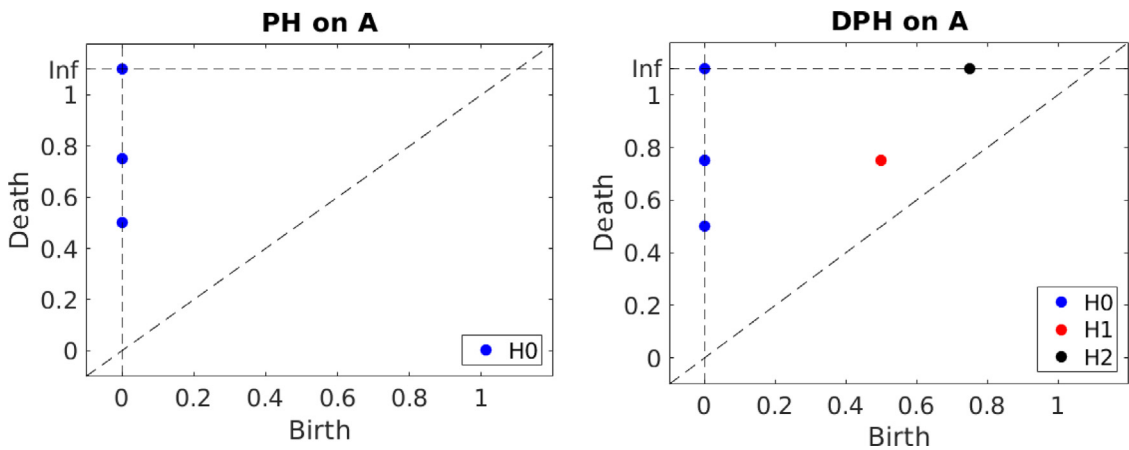


Fig. 4. Left: PH of the matrix A. Right: DPH of the matrix A.

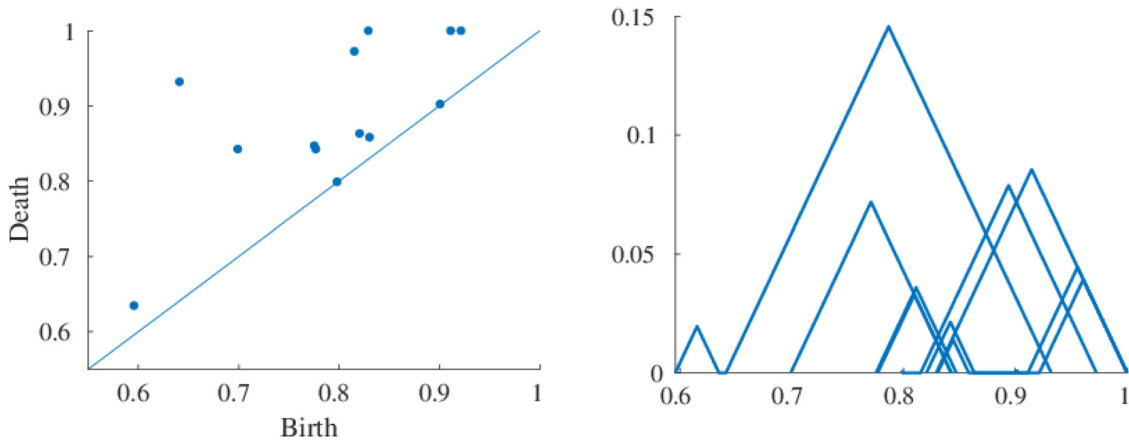


Fig. 5. Left: persistence diagram (sample EEG data). Right: graphical representation of the related persistence landscape.

the set $K \subseteq \mathbb{N}$ to be $K = \{1\}$ to capture the landscape containing information about the most persistent holes; hence, the analysis reported in Section 4 uses the first landscapes $\{\lambda_1^p\}_p$.

3.5.2. Persistence images (PI)

An efficient and stable vector representation of persistent homology is given by the so called *Persistence Images* introduced by Adams et al. (2017). We start with a persistence diagram $D = \{(b_i, d_i)\}_{i \in I}$, a subset of \mathbb{R}^2 as depicted in the example Fig. A.8 and apply the transformation $T : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ sending a birth-death pair (b_i, d_i) to

the birth-persistence pair $(b_i, p_i) := (b_i, d_i - b_i)$. Let $D_k : \mathbb{R}^2 \rightarrow \mathbb{R}$ be the normalized gaussian centered at (b_k, p_k) and with standard deviation σ , i.e.,

$$D_k(x, y) := \frac{1}{2\pi\sigma^2} e^{-[(x-b_k)^2 + (y-p_k)^2]/2\sigma^2}.$$

As set by default in line with Adams et al. (2017), we choose a specific non-negative weighting function $W : \mathbb{R}^2 \rightarrow \mathbb{R}$ that is zero along the horizontal axis, continuous, and piece-wise differentiable, in particular a function depending only on the second coordinate: if p is the maxi-

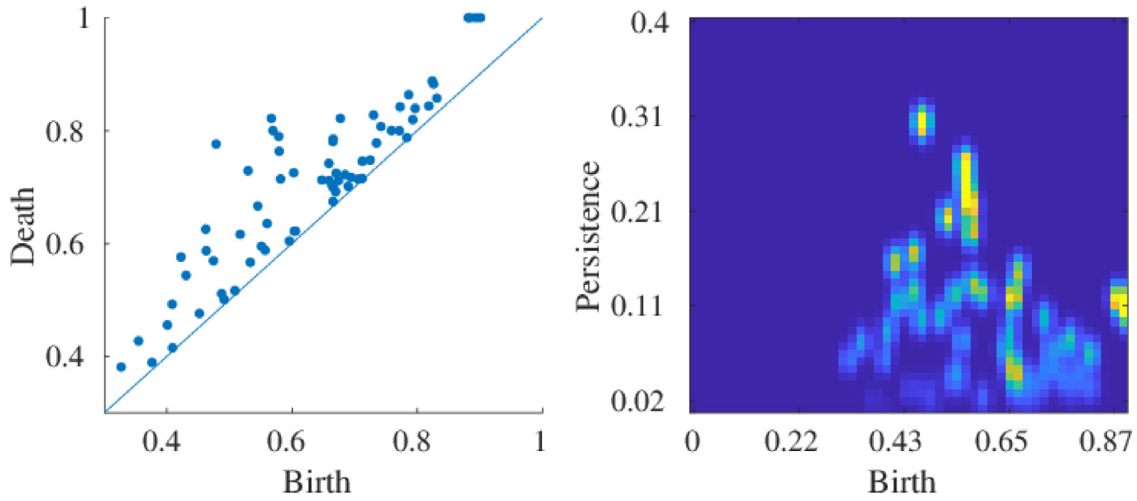


Fig. 6. Left: Persistence diagram of example subject. Right: the associated PI, with domain grid $[0, b_{\max}] \times [0, p_{\max}]$.

mal possible persistence, $W(x, y)$ is set to 0 if $y \leq 0$, 1 if $y \geq p$ and y/p between 0 and 1.

We associate to the persistent diagram D a function $\rho: \mathbb{R}^2 \rightarrow \mathbb{R}$ defined as weighted combination of gaussians:

$$\rho(x, y) := \sum_{i \in I} W(x, y) D_i(x, y).$$

Let b_{\max} and p_{\max} be the maximum birth and maximum persistence across the births and persistences appearing in a diagram, i.e., $b_{\max} := \max_i(b_i)$ and $p_{\max} := \max_i(d_i - b_i)$. After choosing a grid of the domain $[0, b_{\max}] \times [0, p_{\max}]$, we compute the integral $\int \rho(x, y) dx dy$ at each pixel-box of the grid-domain. This gives a discretization of the surface image of $\rho(\mathbb{R}^2)$, hence a matrix of values which is our persistence image. See Fig. 6 for an example of the output.

For a family of persistence diagrams, one can define a *mean PI* of the persistence images associated to the PDs in the family. As each single PI is a matrix of real values, the mean PI of the family is defined as the matrix given by the element-wise average of the matrices. However, this brings a technical issue, as we describe:

Remark 3.2. To construct the PIs, we have used the MATLAB code that accompanies the paper by Adams et al. (2017). Given a Persistence Diagram, one has to set parameters such as the resolution, the probability density and the weighting function. However, as already observed in (Stolz et al., 2018, Section 3.3), in order to compute the mean PI for a family of persistence diagrams, there is an additional pair of values to be chosen: the pair of maximum birth and maximum persistence. We set these values to the maximum across all subjects to achieve unbiased classification results. Note that an alternative group-specific normalization was reported to provide better classification results (Stolz et al., 2018), but in fact provides artificially inflated accuracy due to introducing group-specific bias; a topic that would be covered in more detail in a separate methodological report.

3.5.3. Carlsson coordinates (CC)

Another useful featurization method is given by the so-called *Carlsson Coordinates* (Adcock et al., 2013, Sec. 4.1.2 & 4.2.2). The method is based on an approach from algebraic geometry. For every persistence diagram $D = \{(b_i, d_i)\}_{i \in I}$ in the form of birth-death coordinates, we construct five vector features as follows. Let $p_i = d_i - b_i$ be the life associated to (b_i, d_i) and let p_M be the maximal death time $d_M := \max_i(d_i)$. Following Adcock et al. (2013), we consider the following coordinates:

$$\begin{aligned} f_1(D) &:= \sum_i b_i p_i; & f_2(D) &:= \sum_i p_i (d_M - d_i); \\ f_3(D) &:= \sum_i b_i^2 p_i^4; & f_4(D) &:= \sum_i p_i^4 (d_M - d_i)^2. \end{aligned}$$

Moreover, following (Yesilli et al., 2019, Sec. 4.4), we add a fifth coordinate $f_5(D) := \max_i(p_i)$ taking into account the maximal reached persistence. As suggested in (Adcock et al., 2013, Sec. 4.2.2), we normalize the coordinates f_1, \dots, f_4 by a scale factor $1/N$, where N is the number of the birth-death pairs in the persistence diagram.

3.5.4. Persistent entropy (PE)

Our last feature is *Persistent Entropy*, a topological adaptation of Shannon entropy in the context of persistent homology (Ruocco et al., 2014), i.e., a number that measures the information encoded in the persistent homology groups (at a given dimension).

Let \mathcal{F} be a filtration of a simplicial complex Σ (see Def. Appendix A.4) and $D = \{(b_i, d_i)\}_{i \in I}$ the set of birth-death coordinates associated to PH in dimension k . For every $i \in I$ let $p_i := d_i - b_i$ be the i th persistence.

Definition 3.3. The *persistent entropy* $E(B)$ associated to D is the number $E(D) := -\sum_{i \in I} s_i \log(s_i)$ where $s_i = p_i / S_D$ and $S_D = \sum_{i \in I} p_i$.

3.6. Controlling for inter-individual variability in electrode placement

In real-world scenarios, a substantial variability in the estimated functional or effective connectivity networks between subjects can be due to confounding purely anatomical differences, or spatial misalignment of the observed locations. In an experimental situation with repeated measurements from the same subject, the design may potentially allow suppressing this inter-individual variability by data preprocessing. In particular in iEEG data, a problem of data inconsistency across subjects arises, as every patient has his/her own unique position of electrodes (the implantation) that depends on the particular clinical requirement. Such implantation defines the basic structure of the connections between the recorded brain regions and, as a consequence, iEEG connectivity matrices can not even be directly compared between subjects because of the different data structure and the different sizes.

While this problem is much alleviated in scalp EEG with standardized electrode number and placement, inter-individual differences in brain and head anatomy may still cause some inter-individual variability in EEG connectivity that may obscure intra-individual variability reflecting the brain state dynamics of interest.

Working with the basic topological features rather than with the matrix directly should help to sidestep this problem. However this may not perfectly resolve the issue, as the geometry of the patient's data may still be systematically biased by the individual structure of electrode implantation and anatomy. Therefore, before performing the classification task, we add two steps of data cleaning to the general pipeline.

The first step - used to decrease the influence of implantation on the topology - is directly applied to the connectivity matrices computed by

the procedure described in Section 3.2. The EEG and iEEG data contain epileptic (ictal) and non-epileptic (preictal and interictal) segments as described in Section 3.1.2. In both cases, apart from the ictal and immediately preictal segments, we estimate the connectivity within the clearly interictal segments, which in the following we refer to as the *baseline connectivity*. This baseline connectivity is subtracted from both the preictal and ictal connectivities to provide correction for the confounding inter-individual variability.

The second step in the data cleaning aimed to further control for inter-subject variability is performed after the feature computation and before applying the classification tasks. For each topological feature and patient, we consider the vector consisting of all the realizations of the topological feature across conditions (i.e. all the classified segments for that subject). We compute the mean and remove it from the whole vector of feature values: this gives a final vector with null mean for every patient and feature that has inter-individual biases accounted for, and can be used in the classification tasks.

3.7. Classification

To perform the classification (based on a set of features for each subject), we apply a linear support vector machine (SVM) that constructs a hyperplane in high-dimensional space on a training set of points. This hyperplane is then used for the classification of testing points.

We train a SVM classifier based on the topological features of Section 3.5. We have chosen the SVM as it is the state-of-art method when deep learning can not be applied due to insufficient sample size. Indeed, other models and classifiers could have been used, for example, a random forest approach instead of the SVM. A comparison of these methods employed in classification tasks on persistent homology of FC (albeit in autism rather than schizophrenia) has been recently provided by Rathore et al. (2019), where the methods have shown to achieve similar performances.

Evaluation of the classifier accuracy is made with the leave-two-out procedure on paired samples. In the following, when dealing with the epileptic datasets, by a “data point” we mean (the functional/effective) connectivity matrix characterizing particular brain activity (ictal or inter-ictal); when working with the schizophrenia dataset, “data points” are the single subject connectivities. For each dataset, we train the classifier on the set consisting of all the data points excluding both the tested sample and its matching counterpart. A similar approach has been used, for instance, in the studies by Castro et al. (2011), Mikolas et al. (2018) and Rubin-Falcone et al. (2018). To be more precise, for the epileptic datasets, for every tested connectivity, we exclude its match from the opposite type of activity (ictal/inter-ictal). For the schizophrenia data, for every tested subject, one subject from the opposite group was excluded, in order to keep groups of the same size. For all datasets, presented in the main text, such across-participants leave-two-out procedure was performed. In the supplementary, we also present analysis of EEG data, where classification was performed within-subject in selected subjects with high number of measured seizures, here always a matched pair of ictal and inter-ictal segment from the same recording is removed.

3.8. Software

In the construction of effective connectivity networks we have used the software companion of Barnett and Seth (2013) for the computation of Granger Causality and of Kořenek and Hlinka (2020) for the construction of the mask (see Section 3.2.2). For the calculations of persistent homology groups of undirected networks we have used the python library *Rips*er Tralie et al. (2018). For the computation of persistent homology of the directed graphs, we have used the library *Flagser* by Lütgehetmann et al. (2020) (via the python API reference *pyflagser* of the Giotto-tda Tauzin et al. (2020), version v 0.2.0.), with the default

settings. For computing the PIs, we have used the software that accompanies the paper of Adams et al. (2017); we set the variance of the Gaussians to 0.0001 and keep all the other parameters as by default. The pair maximum birth/maximum persistence is computed across all the instances (see Remark 3.2). For the calculation of PLs, we use our own MATLAB implementation of (Bubenik and Dłotko, 2017, Algorithm 1). For the implementation of support vector machine classifier, we use the MATLAB function *fitsvm*. When otherwise specified, we use the default settings of the declared functions.

4. Results

The aim of the paper is to assess the applicability of (directed) persistent homology methods to brain functional/effective connectivity analysis. We have focused on two prominent neurological applications: epileptic seizure detection (using EEG data, both scalp and intracranial) and classification of schizophrenia/healthy brain states (using fMRI data). We follow the pipeline described in Section 3: for each subject, we first calculate the functional connectivity graph (Section 3.2.1) and the effective connectivity directed graph (Section 3.2.2), hence the associated Vietoris-Rips filtration (Ex. Appendix A.6) and the filtration of directed Flag complex (Def. 2.2), respectively; we then compute the persistent homology groups of the obtained filtrations as described in Section 3.3. We encode the PH information in the form of Persistence Diagrams that we use to compute the associated topological features PL, PI, CC and PE. Here we potentially subtract individual feature mean to control for electrode placement bias as described in Section 3.6 (this is not applicable to the fMRI data where the classification is only between subjects). We finally apply the SVM classifier (see Section 3.7). We compare the obtained classification results with those achieved by using the same SVM classification method applied to Functional Connectivity (FC) (in the undirected case) and to (masked) Effective Connectivity (EC) (in the directed approach). The results are shown in the subsections below. Note that for the results obtained on the iEEG data, a direct comparison with the naive FC/EC approach is not possible, because of different size and position of patient’s implantations.

We proceed by illustrating our results in the group analysis of the fMRI (schizophrenia) data first, of the group analysis of scalp EEG (epilepsy) data afterwards (following the classical pipeline first, and the alternative pipeline later), and we conclude with the analysis of the iEEG (epilepsy) dataset. The persistent homology computations shown in the following subsections refer to calculations in dimension 1 and we refer to Appendix D for the corresponding results in dimension 0.

4.1. First dataset: fMRI - Schizophrenia

In this first analysis, we present the results of a classification using the fMRI time series of 90 healthy subjects and of 90 schizophrenia patients (see Section 3.1.1). We report the results in Table 1.

We observe that the topologically-based classifications are close to random (with accuracy around 0.5) in both the directed and undirected approach. Hence, the schizophrenia patients are not topologically discriminated from the healthy subjects, although the FC-classification is relatively successful (with accuracy 0.79). In this case, the classification task based on the EC matrix itself did not achieve good results (accuracy 0.57, close to the performances of the directed topological features).

4.2. Second dataset: scalp EEG - Epilepsy

Following the same pipeline as in the previous example, we present now a classification analysis concerning a different disease (epilepsy) and a different type of data (EEG). We use for this analysis the EEG recordings described in Section 3.1.

Our analysis of epileptic brain states is divided in the two cases reported below: we first assess a classification by the basic approach as applied to the fMRI data, as described in Section 4.2.1. Subsequently

Table 1

Schizophrenia: group analysis, 90 patients plus 90 healthy controls. SVM classification results using FC, EC and the topological features PL, PI, CC, PE in dimension 1 with undirected approach (PH) or directed approach (DPH). *Acc* stands for *Accuracy*, *Sns* for *Sensitivity* and *Spc* for *Specificity*.

Features		Acc	Sns	Spc
EC	naive	0.57	0.54	0.60
FC	naive	0.79	0.74	0.83
PL	(DPH)	0.44	0.39	0.50
	(PH)	0.48	0.49	0.47
PI	(DPH)	0.49	0.44	0.53
	(PH)	0.56	0.59	0.52
CC	(DPH)	0.34	0.18	0.51
	(PH)	0.61	0.57	0.64
PE	(DPH)	0.54	0.67	0.41
	(PH)	0.56	0.46	0.67

Table 2

Epilepsy: group analysis, 18 subjects. SVM classification results using FC, EC and the topological features PL, PI, CC, PE in dimension 1 with undirected approach (PH) or directed approach (DPH). *Acc* stands for *Accuracy*, *Sns* for *Sensitivity* and *Spc* for *Specificity*.

Features		Acc	Sns	Spc
EC	naive	0.83	0.72	0.94
FC	naive	0.78	0.72	0.83
PL	(DPH)	0.50	0.50	0.50
	(PH)	0.50	0.44	0.56
PI	(DPH)	0.47	0.44	0.50
	(PH)	0.64	0.50	0.78
CC	(DPH)	0.44	0.78	0.11
	(PH)	0.69	0.67	0.72
PE	(DPH)	0.53	0.44	0.61
	(PH)	0.33	0.39	0.28

we report (in Section 4.2.2) the results of the analysis carried on the same dataset but involving the additional cleaning steps described in Section 3.6.

4.2.1. EEG analysis - default approach

In this analysis, we compare the functional and effective connectivity networks associated to the ictal and to the inter-ictal periods. We consider 18 subjects; each subject has both ictal and interictal data segments. For details, see Section 3.1.2. We apply the SVM classifier to the masked EC directed graphs, to the masked FC undirected graphs and to the topological features PL, PI, CC and PE obtained by the PH computation on both the obtained networks. We report the achieved classification values in Table 2.

The results show that the classifications based on the topological features perform again close to random, even though both the raw EC-based and raw FC-based classification are quite effective (with accuracy of 0.83 and 0.78). Therefore, when applied to both the functional and the effective connectivity networks and when compared with the results achieved by the raw matrices, we do not see any advantage in using the topological features.

However, we conjectured that the low performance of the persistent homology approach may be due to its high sensitivity to baseline inter-subject variability in the connectivity structure. Indeed, when applied to the discrimination of ictal and interictal states within a single subject, TDA provides good results, comparable in some cases to those achieved by EC or FC itself, see Appendix C. Therefore, further supported by this successful within-subject application, in the following group analysis we

Table 3

Epilepsy: EEG group analysis (with variation) 18 subjects. SVM classification results using the topological features PL, PI, CC, PE in dimension 1, with undirected approach (PH) or directed approach (DPH). *Acc* stands for *Accuracy*, *Sns* for *Sensitivity* and *Spc* for *Specificity*.

Features		Acc	Sns	Spc
EC	naive	0.44	0.56	0.33
FC	naive	0.89	0.78	1.00
PL	(DPH)	0.58	0.44	0.72
	(PH)	0.75	0.72	0.78
PI	(DPH)	0.86	0.89	0.83
	(PH)	0.83	0.83	0.83
CC	(DPH)	0.61	1.00	0.22
	(PH)	0.92	0.94	0.89
PE	(DPH)	0.89	0.89	0.89
	(PH)	0.06	0.06	0.06

Table 4

Epilepsy: iEEG analysis, 16 subjects. SVM classification results using the topological features PL, PI, CC, PE in dimension 1, with undirected approach (PH) or directed approach (DPH). *Acc* stands for *Accuracy*, *Sns* for *Sensitivity* and *Spc* for *Specificity*.

Features		Acc	Sns	Spc
PL	(DPH)	0.50	0.56	0.44
	(PH)	0.62	0.69	0.56
PI	(DPH)	0.62	0.62	0.62
	(PH)	0.66	0.69	0.62
CC	(DPH)	0.78	0.81	0.75
	(PH)	0.84	0.87	0.81
PE	(DPH)	0.69	0.69	0.69
	(PH)	0.87	0.87	0.87

first attempt to correct for the inter-subject variability, before applying the classifier.

4.2.2. EEG Analysis: Alternative approach

In order to reduce the inter-subject variability of the EEG dataset, in the following group analysis we apply further cleaning steps to the main pipeline, following the methodological details described in Section 3.6. In this analysis we use the same dataset as in Section 4.2.1. The obtained results, in dimension 1, are shown in Table 3; we refer to the Appendix, and in particular to Table D.9, for further results in dimension 0.

We now see, in Table 3, that both DPH and PH achieve on average good performances, with around 70 percent accuracy on average. The directed and undirected method are comparable to FC, whereas EC performs close to random.

The results shown in Table 2 and in Table 3 suggest that the cleaning steps help at increasing the performance of the classification tasks from persistent homology features; on the other hand, the better performances could also be due to some specific properties of the considered data. Finally, we have applied the alternative pipeline to a different dataset, namely to epileptic seizure classification in intracranial EEG data.

4.3. Third dataset: iEEG - Epilepsy

In our last analysis, we use the iEEG dataset described in Section 3.1.2. As in the previous section, in order to reduce the inter-subject variability, we apply the cleaning steps of Section 3.6 before computing the homology features and the classification. The obtained results, are shown in Table 4 and we refer to Table D.9 to additional results in dimension 0.

The results show that good classification performances in the detection of the brain alterations have been achieved, *e.g.*, with values of 0.94 for PE (in dimension 0, see [Table E.11](#)) and 0.87 obtained by PE in dimension 1. These are comparable to the results obtained in the alternative pipeline for scalp EEG data analysis of [Section 4.2.2](#) and the single subject analysis presented in [Appendix C](#) and consolidate the idea that the cleaning step procedures may help in the analysis of brain state detection via TDA. Out here that, in this case, a direct comparison with the performances of the connectivity matrices across subjects via FC/EC is not applicable due to interindividual differences in the electrodes number and placement.

For the sake of completeness, we report in [Table E.11](#) in the appendix also other detailed results of this last analysis, such as employing the full (not-masked) matrices as well and further 0-dimensional features.

5. Discussion

In this work we have investigated the employment of a directed extension of persistent homology to classification tasks of disease-related alterations of brain connectivity from neuroimaging data in the case of two typical brain connectivity alterations, namely an easy task classification between ictal and pre-ictal states, and a more difficult task of classification between schizophrenia patients and healthy controls. The results of the classification of ictal from inter-ictal effective connectivity have provided a proof of principle that the directed persistent homology approach allows extracting features that carry relevant information about the underlying brain state in a real-world setting. However, the failure at the schizophrenia classification task pointed to limitations of the approach. There are several factors that play role here. First, the results of the 'naive' classification using directly the raw connectivity matrices suggest, that it might be that in the case of the fMRI schizophrenia classification task, the estimated effective connectivity (unlike functional connectivity) might have not contained or retained sufficient information to reliably distinguish the patients from the control. In general, effective connectivity aims for richer information concerning the system than functional connectivity (in particular the direction, and in principle the full information about the causal interactions, that can be considered the primary object that gives rise to the secondary statistical dependences captured by the functional connectivity). However, effective connectivity pays for that with higher difficulty of *estimating* it reliably. This may thus also partially explain why even in the seizure classification tasks, the PH approach typically demonstrated similar or better classification performance than the DPH approach, although among the analyses, both situations appear: PH performing better than DPH and vice versa. Indeed, it seems that the most important factor that influences the performances of the classifier is the dataset. We refer to the [Appendix, Section Appendix F](#) for a more detailed discussion of the statistical comparison of the two pipelines, containing permutation-based tests of the classification performances against the null hypothesis of random class assignment, and testing the hypotheses of the EC versus the FC-based classifiers with corresponding feature selection.

As mentioned earlier, the results generally suggest that the naive classifications based on functional or effective connectivity generally perform similarly or better than the topologically-based ones. The better performances of FC/EC on the top of TDA might be due to various reasons, including potentially a lower sensitivity of Persistent Homology to small localized variations of the absolute values of the input adjacency matrices: small variations of the metric properties of relatively small subgraphs, like the distances associated to a subset of nodes and represented by the entries of the adjacency matrix, are loosely distinguished by the associated topological features, but the same variations are better captured by classification methods based on the full raw matrices. This is amplified in the heterogeneous group-studies; in [Appendix C](#), we report an analysis of the EEG epileptic data restricted to the single subjects showing better performances than the group-studies. Driven by the idea that TDA methods would be more effective when tested on more homo-

geneous data, in the second part of our analysis on the EEG datasets, we have introduced further steps in the data processing and cleaning (see [Section 3.6](#)). A theoretically preferable procedure would require to use the covariance matrices instead of the correlations, because satisfying linearity properties and carrying more biological information, but this would involve numerical complications. We observe that, when following this variation of the pipeline, both our computations on the EEG and iEEG data achieve slightly better performances, suggesting that the cleaning procedures may in fact be of importance in the analysis of these data.

Another observation concerns the field of applicability of the persistent homology: for some datasets (as for iEEG data), a classification analysis based on the raw FC/EC matrices can not be undertaken. This is because the analysis assumes fixed number and position of the signal sources, an assumption clearly not fulfilled by the placement of the intracranial electrodes in different subjects. Although under some circumstances homogenous resampling might be attempted (such as averaging of signals within anatomical regions), investigations like those shown on the iEEG data might be the natural field of application of TDA within neuroimaging. We note here that other measures such as raw FC/EC with some dimension reduction or indeed other methods not using connectivity at all are likely to provide similar or even higher classification performance than that achieved by the current combinations of connectivity and persistent homology features, as the seizure network dynamics are indeed commonly quite different from the inter-ictal period. We refer readers interested in this particular classification task to a recent review by [Siddiqui et al. \(2020\)](#), however noting that for a range of the reviewed studies, only the best performing classifier is reported, and the accuracies might therefore be generally inflated.

Due to richness of the available methodologies, we necessarily had to limit our investigation by making some specific choices. Among the general advantages of PH over more classical methods, PH does not depend upon the choice of a threshold. However, it does depend on the choice of the filtrations used and simpler graph filtrations (see, *e.g.*, [Lee et al. \(2011a\)](#)) may also be employed as input of PH instead of the more general Rips/cliue filtration. Although both the graph and Rips/cliue filtrations lead to the same features in dimension 0 (as the number of connected components of the graphs and of the associated simplicial complexes is the same), this is not true in higher dimensions. In fact, the n -dimensional features obtained by using the graph filtration are 0 for every $n \geq 2$, whereas, due to the presence of higher-order simplices, when using the Rips/cliue filtration non-trivial n -dimensional holes may arise. In this work we have used the more general cliue filtration, that in principle provides additional information due to the presence of higher-order simplices; however this does not necessarily mean that the classifiers will be able to use this information efficiently; however, we have not provided a comparison with other filtrations. Among other specific choices, for instance, many more featurizations of persistent homology are available in the literature, and we thus decided to focus only on some of the more common. Among other examples, kernel methods are also prominent, and we report some kernel-based classifications in [Appendix D](#). The results reported in [Section 4](#) and in the [Appendix](#), refer to computations using features in dimension 0 and 1; we have also computed the persistent homology groups in higher dimensions, but, due to the size of our datasets, the Betti numbers were typically zero (and always zero from dimension 3 higher), so we did not include them in the report. It is an open question whether, for larger datasets, the higher Betti numbers can be effectively used for reaching better classifications.

Moreover, there is the choice of potentially discarding the weak (potentially noise-related) links from consideration. Among prior studies analysing functional connectivity networks through persistent homology methods, some set statistically insignificant correlations (with p -values greater than 0.05) to 0 (see [Merelli et al. \(2016\)](#); [Stolz et al. \(2018\)](#)) and others do not (see [Lee et al. \(2011b, 2012\)](#); [Rathore et al. \(2019\)](#)). The results reported in this work refer to connec-

tivity networks with links surviving significance threshold of $\alpha = 0.05$, but either approach can be equivalently undertaken; for sake of completeness, we report in Table D.7 the results obtained by using the full EC or FC matrix (from which we compute the distance matrix by the transformation $1 - |\cdot|$) and a comparison of the performances, showing that the performance does not substantially depend on the choice of these masks. For consistency with the undirected approach, in our study of the effective connectivity networks we introduce a similar masking approach: we first build the (pair-wise) Granger Causality graph associated to the brain data and then we set to 0 the connections in the network marked as dependent by applying the FACDA algorithm (Kofenek and Hlinka, 2020). Table D.7 shows that also in the case of the effective network, masking does not change the results substantially.

Another specific choice is the use of the SVM rather than other classifier. As mentioned in the Methods section, SVM is a relatively standard method of choice for similar data situations, used also previously in similar context Rathore et al. (2019); Stolz et al. (2018), although a plethora of alternatives exist. As with most classification tools, the performance of SVM might be adversely affected by high dimensionality of input data. However, except for PL and, potentially, PI, all the other TDA features that we have used are low-dimensional, in particular in comparison to the number of features in the naive approach that uses the full connectivity matrix as its input.

In order to study effective connectivity networks from a topological point of view, there are two approaches: one can either construct an undirected graph from the directed effective connectivity network (hence losing some information) and then apply the usual PH pipeline or one can use a suitable directed version of PH. To the best of our knowledge, this is the first attempt at employing the latter, directed PH approach to analyse effective connectivity. However, this approach has been foreseen in a recent perspective paper by Lee (2019), suggesting that causal connectivity methods are to be examined and applied for classification of diseases.

To summarise, we have shown that the (directed) persistent homology approach can be applied to brain functional/effective connectivity, obtaining in some cases good performances. However we suggest that a due caution in the use of TDA and of these featurizations should be taken, in line with very recent critical results on TDA based on functional connectivity in the literature Ellis et al. (2019); Rathore et al. (2019).

6. Conclusion

In this work we have shown that in principle the directed extension of persistent homology can be successfully applied to classification of disease-related alterations of brain connectivity from neuroimaging data. However, cautious direct comparison to the performance of standard (undirected) persistent homology features as well as to the use of the raw effective (or functional) connectivity matrices for the same classification task suggests that there are important challenges to be overcome to fully utilize the theoretical promise of the topological data analysis. We conclude that tackling the identified real-world neuroimaging challenges such as topographical versus topological specificity of the brain alterations, inter-individual biases and homogeneous versus heterogeneous sampling situations, would be necessary in order to firmly establish TDA also in the particular field of neuroimaging as a pragmatic tool.

Credit authorship contribution statement

Luigi Caputi: Conceptualization, Methodology, Software, Formal analysis, Investigation, Writing - original draft, Writing - review & editing, Visualization. **Anna Pidnebesna:** Methodology, Software, Validation, Formal analysis, Investigation, Data curation, Writing - review & editing, Visualization. **Jaroslav Hlinka:** Conceptualization, Methodology, Formal analysis, Resources, Writing - review & editing, Supervision, Funding acquisition.

Acknowledgment

The authors thank Jakub Kopal for sharing the preprocessed fMRI time series and Barbora Bučková for sharing scripts for classification pipeline.

Appendix A. Background: Topological Data Analysis

Topological Data Analysis is a recent and active research field of applied algebraic topology, whose main goal is to study the structure of datasets by means of topological frameworks. Among the tools of TDA maybe the most common is nowadays Persistent Homology, whose definition is based on the mathematical ideas of *simplicial complexes*, *filtrations* and *homology*. In this section we present a brief overview on this subject, we describe these concepts in the sections Section A.1 and Section A.2 below (referring to Otter et al. (2017); Zomorodian and Carlsson (2005) for a mathematical introduction and to Horak et al. (2009); Petri et al. (2013) for further studies of Persistent Homology with applications to complex networks).

A1. Simplicial complexes

Simplicial complexes are combinatorial mathematical objects that generalize (in higher dimensions) the classical definition of graphs. Roughly speaking, a simplicial complex is a space built by assembling together points, edges, triangles, tetrahedra and, more generally, high dimensional polytopes (in topology called simplices). We refer to Munkres (1984) for more details.

Definition A.1. Let V be a set. An (abstract) simplicial complex on V is a collection $\Sigma = \{\sigma_\alpha\}$ of finite non-empty subsets σ_α of V such that:

- (i) for each $v \in V$, $\{v\}$ belongs to Σ ;
- (ii) if $\sigma \in \Sigma$ and $\tau \subseteq \sigma$, then $\tau \in \Sigma$.

The elements σ_α of a simplicial complex Σ are called *simplices* of Σ and if σ_β is a subset of a simplex σ_α , then we say that σ_β is a *face* of σ_α . The *dimension* of a simplex σ is one less its cardinality, i.e., if $\sigma = \{v_0, \dots, v_k\}$ then $\dim \sigma = k$, and we say that σ is a k -simplex. The 0-simplices $\{v\} \in \Sigma$ are called the *vertices* of the simplicial complex Σ . We will write $V(\Sigma)$ for the set of vertices of Σ ; with abuse of notation, we will make no distinction between the element $v \in V$ and the vertex $\{v\} \in V(\Sigma)$.

In this work we consider only finite simplicial complexes, so that the family $\Sigma = \{\sigma_\alpha\}$ is finite. The dimension of the simplicial complex Σ is the maximal dimension across all its simplices.

Observation. Definition Appendix A.1 is rather abstract. However, one can always geometrically realize an (abstract) simplicial complex as a geometric subspace of \mathbb{R}^n . A k -simplex $\sigma = \{p_0, \dots, p_k\}$ can be geometrically depicted as the convex hull of $k + 1$ geometrically independent points p_0, \dots, p_k in \mathbb{R}^{k+1} . For example, a 0-simplex is a point, a 1-simplex is depicted as a segment (a convex hull of its two endpoints, which are again 0-simplices). Likewise, a 2-simplex is geometrically understood as a triangle, and so on. Roughly, for geometrically depicting a simplicial complex, one realizes each of its simplex as a subspace of \mathbb{R}^n and then glues them to each other along common faces. See Fig. A.7 for an example of geometric realizations of different simplicial complexes (there denoted by $\Sigma_0, \dots, \Sigma_5$).

Example A.2. Recall that a graph G is a pair $G = (V, E)$ given by a set V of vertices and a set $E \subseteq V \times V$ of edges, where an edge is identified by the couple (v, w) of its distinct endpoints (so we do not allow multiple edges or loops). Hence, a graph is a simplicial complex whose vertex set is V and the 1-simplices are the edges of G .

If Σ and Σ' are simplicial complexes, a simplicial map $\Sigma \rightarrow \Sigma'$ between them is a function $f : V(\Sigma) \rightarrow V(\Sigma')$ sending vertices of Σ to vertices of Σ' , with the property that, if the vertices v_0, \dots, v_n span a simplex

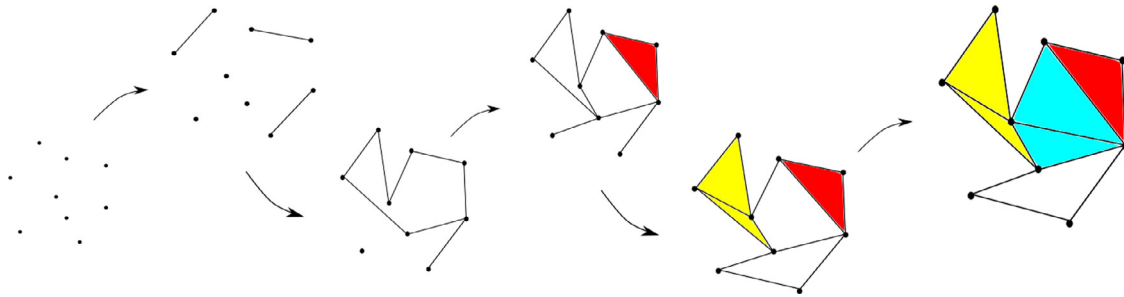


Fig. A1. Example of a filtration of simplicial complexes. At each step we add some edges, and when we find a 3-clique we fill it in with a (coloured) 2-simplex.

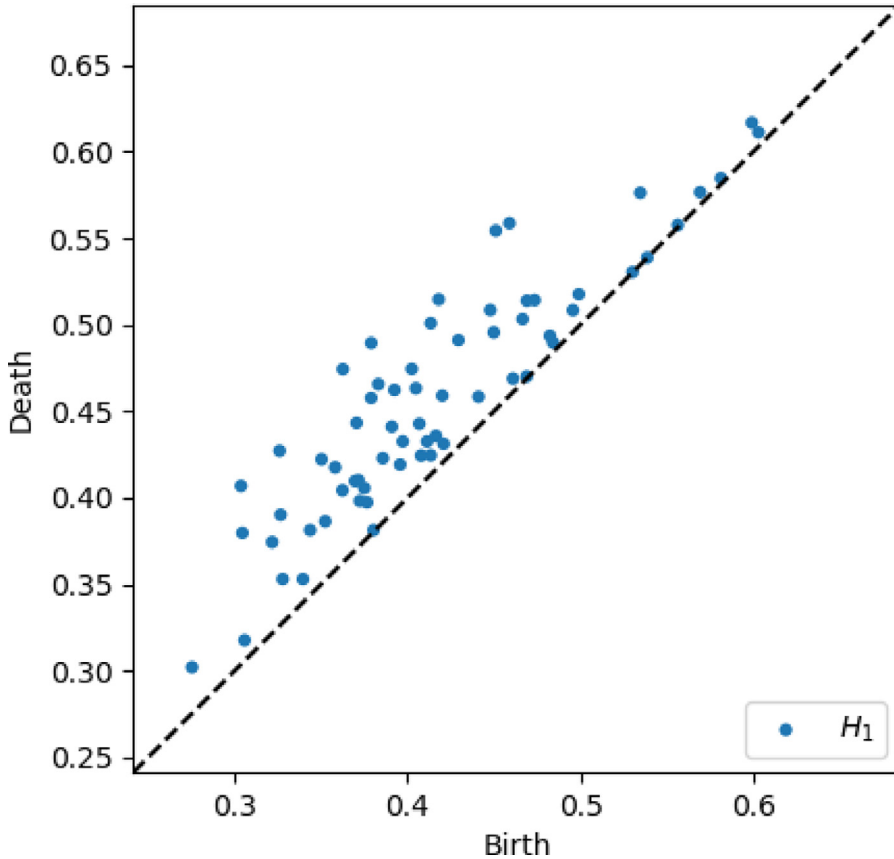


Fig. A2. Persistence diagram associated to a filtration of simplicial complexes. Each dot represents the birth-death of a k -dimensional hole. The bottom right triangular segment is empty because the death of a persistent homology class can not be lower than its birth.

$\sigma = \{v_0, \dots, v_n\}$ of Σ , then also the vertices $f(v_0), \dots, f(v_n)$ span a simplex of Σ' ; observe that the image of σ can be a simplex of Σ' of smaller dimension. Examples of simplicial maps arise from inclusion. For example, if Σ is contained in Σ' (as subfamily) then the inclusion $\Sigma \subseteq \Sigma'$ defines a map of simplicial complexes (sending the simplices of Σ to the same simplices seen in Σ'). In this case we say that Σ is a *subcomplex* of Σ' . Having maps between simplicial complexes preserving the structures is an important ingredient in developing the theory of filtrations and persistent homology (see the sections below).

Example A.3 (Clique complex). A graph G is a 1-dimensional simplicial complex (as its simplices are at most of dimension 1); recall that an n -clique is a complete subgraph of G on n vertices. The *clique complex* \hat{G} associated to G is the simplicial complex constructed as follows:

1. its 0- and 1-simplices are the vertices and edges of G ;
2. for every $(n + 1)$ -clique with vertices $\{v_0, \dots, v_n\}$ in G we add the n -simplex $\sigma_\alpha = \{v_0, \dots, v_n\}$ on the same vertices.

For example, if G is a complete graph on 3 vertices, then its clique complex is nothing but a simplicial complex consisting of a 2-simplex

together with all its faces. Observe that the dimension of the clique complex depends on the dimension of the maximal cliques of G . The embedding of G into \hat{G} sending the vertices and edges of G to the respective 0 and 1-simplices of \hat{G} is a map of simplicial complexes.

We are ready to recall the definition of filtration of a simplicial complex Σ :

Definition A.4. A *filtration* \mathcal{F} (of a simplicial complex Σ) is an indexed sequence of subcomplexes of Σ

$$\mathcal{F} : \Sigma_0 \subseteq \dots \subseteq \Sigma_N = \Sigma$$

such that for every $i \leq j$, the simplicial complex Σ_i is contained in Σ_j .

We will only consider finite filtrations of simplicial complexes (so that the chain of simplicial complexes in the definition is finite). The set of indices in the definition of a filtration is usually given by positive natural numbers, but it can also be a more general set, like the positive real numbers (see, e.g., [Example Appendix A.6](#)). A figurative illustration of a filtration (with clique complexes) is given in [Fig. A.7](#).

We conclude the subsection with two prominent examples of filtrations: the *weight rank clique filtration* (WRCF) Petri et al. (2013) associated to a weighted graph and the *VietorisRips filtration* associated to any finite subset of a metric space.

Example A.5 (WRCF). Let G be a weighted graph and let $(w_i)_{0 \leq i \leq N}$ be the increasing sequence of all the weights of G (i.e., w_0 is the minimal weight, w_N the maximal and $w_i < w_{i+1}$). For the i th weight w_i construct the simplicial complex Σ_i as follows:

1. consider the subgraph $G_i \subseteq G$ with the same vertices of G and only those edges of G with corresponding weight w such that $w \leq w_i$;
2. $\Sigma_i := \hat{G}_i$ is the clique complex of G_i (see Ex. Appendix A.3).

It is easy to see that Σ_i is contained in Σ_{i+1} and that the sequence $\emptyset \subseteq \Sigma_0 \subseteq \dots \subseteq \Sigma_N = \Sigma$ is a filtration. Analogously, by sorting the weights in the decreasing order, one obtains another filtration, where the graph G_i is defined by keeping only the edges of G with weight $w \geq w_i$.

Example A.6 (VietorisRips filtration). Let (X, d) be a metric space and S be a finite subset of X , endowed with the induced metric. In applications, the set S can be a set of data-points in \mathbb{R}^n or a set of measurements with a notion of distances between them. The *VietorisRips complex of S at scale $t \geq 0$* is the simplicial complex:

$$VR_t(S) = \{\emptyset \neq \sigma \subseteq S \mid \text{diam}(\sigma) \leq t\}$$

where the $\text{diam}(\sigma) := \sup\{d(x, y) \mid x, y \in \sigma\}$ is the diameter of σ (we point out here that different definitions, considering a maximal distance of $2t$ instead of t can be found in the literature; we follow the definition used in Tralie et al. (2018) to be consistent with Section 3.8). Observe that every such σ represents a simplex of dimension $|\sigma| - 1$. If $t \leq t'$ then $VR_t(S) \subseteq VR_{t'}(S)$ and the resulting filtration $VR(S)$ of the simplicial complexes $VR_t(S)$, with t varying in $\mathbb{R}_{\geq 0}$, is called the *VietorisRips filtration of S* .

Remark A.7. Observe that the two constructions in Ex. Appendix A.5 and in Ex. Appendix A.6 are equivalent, if one considers metric spaces (X, d) associated to weighted graphs G (with respect to an increasing ordering of the weights), with underlying space X given by the nodes of the graph G and distance $d(x, y)$ between x and y defined as the weight corresponding to the edge (x, y) .

A.2. (Persistent) homology and persistence diagrams

Persistent Homology (PH) Edelsbrunner et al. (2002); Zomorodian and Carlsson (2005) is one of the main tools in TDA and one of the first applications of Algebraic Topology to Data Analysis, introduced in Carlsson (2009). Persistent Homology is a topological invariant whose definition is based on the mathematical concept of homology groups.

It is robust to noise and stable with respect to small perturbations of the inputs (Cohen-Steiner et al., 2007, Section 3.1); moreover, calculations can be easily achieved by using many software options (in MATLAB, Python, C++, R, Java) Otter et al. (2017). Its constituents, the homology groups, are classical algebraic invariants of topological spaces. The number of connected components, loops (1-dimensional holes), voids (2-dimensional holes), tunnels and higher dimensional holes are examples of such invariants. A technical comprehensive definition of homology groups is outside of the scope of the present work. In the following, we will only sketch some ideas, but the interested reader can find more details in Hatcher (2000); Munkres (1984). We conclude the section with the definition of persistent homology groups and persistence diagrams.

Let k be a natural number. The k th *homology group* $H_k(\Sigma)$ associated to a simplicial complex Σ (see Def. Appendix A.1) can be thought of as the set (group, module or vector space) of k -dimensional holes in Σ ; the elements of $H_k(\Sigma)$ are called *homology classes*. To be more precise, let \mathbb{F}_2 be the field with two elements (i.e., 0 and 1, with sum and product

inherited from the usual sum and product of real numbers, reduced mod 2) and let $C_p(\Sigma)$ be the free \mathbb{F}_2 -vector space whose basis consists of the set of p -simplices of Σ : elements of $C_p(\Sigma)$ are formal combinations $a_0\sigma_0 + \dots + a_n\sigma_n$ of p -simplices $\sigma_0, \dots, \sigma_n$ of Σ with coefficients a_i in \mathbb{F}_2 . For every $p \geq 1$ we define the map

$$d_p : C_p(\Sigma) \rightarrow C_{p-1}(\Sigma), \quad d_p(\sigma) := \sum_{\tau \subseteq \sigma, \tau \in C_{p-1}(\Sigma)} \tau$$

by sending a p -simplex σ to a formal sum of all its $(p-1)$ -faces. The map d_0 is defined as the zero map. As the composition $d_p \circ d_{p+1} = 0$ is the zero map, the image of d_{p+1} in $C_p(\Sigma)$ is contained in the kernel of d_p ; these are both \mathbb{F}_2 -vector spaces and the quotient $\ker(d_k)/\text{Im}(d_{k+1})$ of vector spaces is well-defined. We can now formally define the homology groups of the simplicial complex Σ (with \mathbb{F}_2 -coefficients):

Definition A.8. Let Σ be a simplicial complex and $k \geq 0$ a natural number. The k th *homology group* $H_k(\Sigma)$ of Σ

$$H_k(\Sigma) := \ker(d_k)/\text{Im}(d_{k+1})$$

is defined as the quotient of the kernel of the map d_k with the image of d_{k+1} .

The group $H_k(\Sigma)$ is in fact a \mathbb{F}_2 -vector space and its dimension (over \mathbb{F}_2) called the k th *Betti number* $\beta_k(\Sigma)$ of Σ . More general definitions are possible, with arbitrary fields \mathbb{F} or with the integers \mathbb{Z} as coefficients. We refer to Hatcher (2000); Munkres (1984) for more detailed descriptions of homology groups.

The number of k -holes of (the geometric realization of) a simplicial complex Σ corresponds to the k th Betti number $\beta_k(\Sigma)$ associated this way to Σ . In fact, Betti numbers give (a measure of) the number of “independent” holes or loops in Σ , i.e., holes that can not be shrunk continuously to a point. The 0th Betti number β_0 represents the number of connected components, the 1-st Betti number β_1 is the number of independent loops and the 2-nd Betti number β_2 gives the number of 2-dimensional voids.

Example A.9. If X is a point, then X has a single connected component and no loops or higher dimensional holes. Its Betti numbers $\beta_n(X)$ are 0 for every $n > 0$, except for $\beta_0(X)$ which is 1. The 0th Betti number of two points is 2 (2 connected components) and 0 in higher dimensions. On the other hand, if X is a 2-dimensional sphere, then X has a 2-dimensional hole, represented by the void internal to the sphere. We get $\beta_0(X) = \beta_2(X) = 1$ and $\beta_n(X) = 0$ otherwise. Observe that $\beta_1(X)$ is also 0 as every loop on the sphere can be continuously shrunk to one point.

One of the most useful theoretical properties of homology is given by its functoriality: for every simplicial map $f : X \rightarrow Y$ of simplicial complexes and for every k in \mathbb{N} , we also get a linear map $H_k(f) : H_k(X) \rightarrow H_k(Y)$:

$$f : X \rightarrow Y \quad \mapsto \quad \{H_k(f) : H_k(X) \rightarrow H_k(Y)\}_{k \in \mathbb{N}} \tag{A.1}$$

satisfying some additional compatibility properties; we avoid to spell them out here and we refer to Hatcher (2000) for the details. This turns out to be of great importance when dealing with more than one single simplicial complex.

As we have seen, the k th homology group $H_k(\Sigma)$ is associated to a single simplicial complex Σ . Likewise, persistent homology groups are defined for a filtration of simplicial complexes, as we now explain. By the functoriality described in Eq. A.1, if $F : \emptyset \subseteq \Sigma_0 \subseteq \dots \subseteq \Sigma_N = \Sigma$ is a filtration of simplicial complexes (see Def. Appendix A.4), then, for every k , we get a sequence $0 \rightarrow H_k(\Sigma_0) \rightarrow \dots \rightarrow H_k(\Sigma_N) = H_k(\Sigma)$ of homology groups and morphisms between them. Let k be fixed; then, for every pair of indices i and j between 0 and N , the composition of maps of simplicial complexes $\Sigma_i \rightarrow \dots \rightarrow \Sigma_j$ induces a map $f_k^{i,j} : H_k(\Sigma_i) \rightarrow H_k(\Sigma_j)$ between the homology groups of Σ_i and Σ_j .

Definition A.10. Edelsbrunner et al. (2002) For $i < j$, the k th *persistent homology group* $H_k^{i,j}(F)$ of the filtration $F : \emptyset \subseteq \Sigma_0 \subseteq \dots \subseteq \Sigma_N = \Sigma$ is the

image of the induced map $f_k^{i,j} : H_k(\Sigma_i) \rightarrow H_k(\Sigma_j)$. The k th persistent Betti numbers $\beta_k^{i,j}(F)$ are the ranks of the persistent homology groups $H_k^{i,j}(F)$.

The intuition behind the definition of persistent homology groups is that homology classes (the k -dimensional holes) can appear and disappear at different stages of the filtration steps, due to changes in the topology of the simplicial complexes across the filtration. A persistence homology class (i.e., an element) of $H_k^{i,j}(F)$ represents a k -dimensional hole of Σ_i surviving across all the filtration steps between the index i and index j . The filtration indices at which a k -dimensional hole first appears and first disappears are called its *birth* and *death*. Persistent homology makes (mathematically) clear the idea of appearing and disappearing of high dimensional holes; it keeps track of all these topological changes with the persistence Betti numbers.

We can visualize the information provided by persistent homology groups and persistent Betti numbers by constructing a planar diagram called *Persistence Diagram* (PD), stably with respect to the input data Cohen-Steiner et al. (2007). A persistent homology class born at time b and that died at time d is represented in the diagram by the pair (b, d) (see Fig. A.8). Classes that persist until the end are usually depicted in the persistent diagram as dying at infinity (or at any other index bigger than the maximum filtration one). Note that there are no points below the diagonal because the death value of a hole is always bigger than its birth value.

In a persistence diagram it is then visually clear the measure of how long a homology class persists: the farther from the diagonal, the higher the persistence. In the standard interpretations, short-living classes (so, with representing dots close to the diagonal) are interpreted as noise, and only long-persistence classes represent important features of the filtration.

Appendix B. EEG data: Subject selection

In our analysis of epileptic brain states we have used the data available at <https://physionet.org/content/chbmit/1.0.0/>, from the publicly open dataset *Physionet*. We have analysed EEG recordings belonging to 18 selected subjects: chb01-05, chb07, chb09, chb10, chb12, chb14, chb17-24 except chb18_01, chb19_01.

For the variation of the main pipeline, we have used the following files: chb01_03, chb02_16+, chb03_01, chb04_05, chb05_06, chb07_12, chb09_06, chb10_12, chb12_06, chb14_11, chb17a_03, chb18_29, chb19_28, chb20_13, chb21_19, chb22_20, chb23_06, chb24_04. For those, we always use the first seizure (the first 30 s of it) and two 30-second segments preceding it: [-3m:-2m 30s] for baseline connectivity, and [-1min:-30sec] for preictal connectivity estimation. Detailed procedure is described in Section 3.1.2, in Intracranial EEG data description.

Appendix C. EEG data: single subject analysis

In this supplementary, we assess the efficiency of PH and DPH on single subjects brain state alteration classification, in contrast to the group analysis reported in the main paper in Section 4.2.1.

In the single-subject analysis, we consider the data collected from the patients CHB06, CHB15 and CHB24 (from the same database, we have chosen the patients with higher number of seizures).

For every patient, we consider segments with epileptic seizures and interictal data of corresponding length: 10 ictal plus 10 interictal recordings for the patient CHB06, 20 plus 20 for the patient CHB15 and 10 plus 10 for the patient CHB24. We then apply the same pipeline as in the group analysis.

From the accuracy results, shown in Table C.5, we observe that the classification performances associated to single patients, based on topological features out of the directed and undirected approach, achieve better results (on average about 0.75, 0.60 and 0.70 for CHB06, CHB15 and CHB24, respectively) than in the inter-subject analysis.

Table C.5

Epilepsy: single subject analysis. SVM classification of ictal and interictal data corresponding to 3 patients: CHB06, CHB15, CHB24. For each patient, the classification uses masked FC, masked EC or the topological features PL, PI, CC and PE with undirected approach (PH) and with directed approach (DPH). Acc stands for Accuracy, Sns for Sensitivity and Spc for Specificity.

Patient	Features	Method	Acc	Sns	Spc	
CHB06	EC	(DPH)	0.95	0.90	1.00	
		(PH)	1.00	1.00	1.00	
	FC	(DPH)	0.65	0.80	0.50	
		(PH)	0.60	0.70	0.50	
	PL	(DPH)	0.75	0.60	0.90	
		(PH)	0.75	0.80	0.70	
	PI	(DPH)	0.80	1.00	0.60	
		(PH)	0.85	0.90	0.80	
	CC	(DPH)	0.75	0.60	0.90	
		(PH)	0.95	0.90	1.00	
	CHB15	EC	(DPH)	0.78	0.80	0.75
			(PH)	0.93	0.90	0.95
FC		(DPH)	0.48	0.40	0.55	
		(PH)	0.50	0.50	0.50	
PL		(DPH)	0.60	0.55	0.65	
		(PH)	0.68	0.70	0.65	
PI		(DPH)	0.63	0.70	0.55	
		(PH)	0.45	0.50	0.40	
CC		(DPH)	0.48	0.05	0.90	
		(PH)	0.30	0.30	0.30	
CHB24		EC	(DPH)	0.70	0.90	0.50
			(PH)	0.90	0.90	0.90
	FC	(DPH)	0.65	0.60	0.70	
		(PH)	0.85	0.80	0.90	
	PL	(DPH)	0.65	0.80	0.50	
		(PH)	0.80	0.70	0.90	
	PI	(DPH)	0.75	0.80	0.70	
		(PH)	0.60	0.50	0.70	
	CC	(DPH)	0.55	0.40	0.70	
		(PH)	0.95	0.90	1.00	

Topological features may not achieve great performances in heterogeneous group-wise classification settings (compare the results with those in Table 2), but they are more informative when applied to more specific datasets.

Appendix D. Supplementary results

In a parallel exploratory analysis to that presented in the main text, we have applied small variations of the pipeline described in Section 3 and/or computed different features. In this supplementary we show the results obtained following the main pipeline applied on the fMRI and EEG datasets, using either kernel methods (see Appendix D.1), the full matrices (as opposed to the masked ones, see Appendix D.2, or the 0-dimensional features (as opposed to the 1-dimensional features, see Appendix D.3).

D1. Kernel methods

Our main analysis has been based on Persistence Images, Persistence Landscapes, Persistence Images, Carlsson Coordinates and Persistent Entropy, but many other topological features to analyse datasets are available. Among others, kernel-based methods provide some flexible choices. As there is not a standard persistence kernel approach, we have decided to analyse this feature separately. We follow the kernel-based method introduced by Reininghaus et al. (2015) and apply the SVM classification on the (masked) FC and EC networks, following the pipeline in Section 3. We report the results in Table D.6, from which we can see that they are comparable to those achieved by the other features reported in the main text of the work.

Table D.6

SVM kernel-based classification of functional and effective connectivity networks (both masked) using undirected approach (PH) or directed approach (DPH) (in dimension 1), respectively. Epil denotes the EEG epilepsy data for the group analysis, CHB06, CHB15 and CHB24 denote single epileptic patients and schizo denotes the fMRI schizophrenia data. Acc stands for Accuracy, Sns for Sensitivity and Spc for Specificity.

Data	Method	Acc	Sns	Spc
SchizoEC	(DPH)	0.51	0.57	0.46
	(PH)	0.62	0.67	0.57
Epil	(DPH)	0.53	0.39	0.67
	(PH)	0.61	0.67	0.56
CHB06EC	(DPH)	0.80	0.90	0.70
	(PH)	0.85	0.80	0.90
CHB15EC	(DPH)	0.43	0.50	0.35
	(PH)	0.58	0.70	0.45
CHB24EC	(DPH)	0.55	0.60	0.50
	(PH)	0.75	0.70	0.80

Table D.7

Full FC vs masked FC; full EC vs masked EC. Epil denotes the EEG epilepsy data for the group analysis, CHB06, CHB15 and CHB24 denote single epileptic patients and schizo denotes the fMRI schizophrenia data. The FC/EC masks are discussed in Section 3.2. Features refer to the topological features of Section 3.5; naive denotes the classification method applied to the original matrices.

Data	Feature	Full-FC	Masked-FC	Full-EC	Masked-EC
SchizoCC	PE	0.59	0.61	0.59	0.34
	PI	0.58	0.56	0.55	0.54
	PL	0.52	0.56	0.48	0.49
	naive	0.53	0.48	0.59	0.44
	naive	0.80	0.79	0.54	0.57
Epil	CC	0.56	0.69	0.58	0.44
	PE	0.56	0.33	0.61	0.53
	PI	0.58	0.64	0.33	0.47
	PL	0.36	0.50	0.31	0.50
	naive	0.81	0.78	0.89	0.83
CHB06CC	CC	1.00	0.85	0.70	0.80
	PE	0.80	0.95	0.60	0.75
	PI	1.00	0.75	0.65	0.75
	PL	0.95	0.60	0.60	0.65
	naive	1.00	1.00	1.00	0.95
CHB15CC	CC	0.28	0.45	0.58	0.63
	PE	0.48	0.30	0.48	0.48
	PI	0.50	0.68	0.63	0.60
	PL	0.73	0.50	0.65	0.48
	naive	0.93	0.93	0.95	0.78
CHB24CC	CC	0.90	0.60	0.65	0.75
	PE	1.00	0.95	0.60	0.55
	PI	0.75	0.80	0.40	0.65
	PL	0.95	0.85	0.70	0.65
	naive	0.80	0.90	0.85	0.70

Table D.8

SVM classification with 0-dimensional topological features, for the masked FC and EC. Epil denotes the EEG epilepsy data for the group analysis, CHB06, CHB15 and CHB24 denote single epileptic patients and schizo denotes the fMRI schizophrenia data. Features (CC) and (PE) refer to the topological features of Section 3.5, (kernel), (AUC) and (IPF) are additional features discussed in the text.

Data	Features	Method	Acc	Sns	Spc
SchizoCC		EC (DPH)	0.52	0.30	0.73
		FC (PH)	0.49	0.51	0.47
PE		EC (DPH)	0.49	0.19	0.80
		FC (PH)	0.49	0.62	0.37
kernel		EC (DPH)	0.53	0.30	0.77
		FC (PH)	0.53	0.59	0.47
AUC		EC (DPH)	0.49	0.70	0.28
		FC (PH)	0.55	0.56	0.54
IPF		EC (DPH)	0.45	0.62	0.28
		FC (PH)	0.38	0.19	0.57
Epil	CC	EC (DPH)	0.47	0.50	0.44
		FC (PH)	0.53	0.61	0.44
PE		EC (DPH)	0.33	0.00	0.67
		FC (PH)	0.53	0.72	0.33
kernel		EC (DPH)	0.42	0.50	0.33
		FC (PH)	0.53	0.67	0.49
AUC		EC (DPH)	0.44	0.50	0.39
		FC (PH)	0.56	0.61	0.50
IPF		EC (DPH)	0.53	0.78	0.28
		FC (PH)	0.56	0.50	0.51
CHB06CC		EC (DPH)	0.75	0.70	0.80
		FC (PH)	1.00	1.00	1.00
PE		EC (DPH)	0.70	0.80	0.60
		FC (PH)	0.95	1.00	0.90
kernel		EC (DPH)	0.75	0.80	0.70
		FC (PH)	1.00	1.00	1.00
AUC		EC (DPH)	0.65	0.70	0.60
		FC (PH)	1.00	1.00	1.00
IPF		EC (DPH)	0.80	0.90	0.70
		FC (PH)	1.00	1.00	1.00
CHB15CC		EC (DPH)	0.30	0.20	0.40
		FC (PH)	0.58	0.55	0.60
PE		EC (DPH)	0.60	0.55	0.65
		FC (PH)	0.58	0.50	0.65
kernel		EC (DPH)	0.50	0.40	0.60
		FC (PH)	0.38	0.55	0.20
AUC		EC (DPH)	0.58	0.45	0.70
		FC (PH)	0.58	0.75	0.40
IPF		EC (DPH)	0.63	0.55	0.70
		FC (PH)	0.50	0.20	0.80
CHB24CC		EC (DPH)	0.70	0.70	0.70
		FC (PH)	0.95	1.00	0.90
PE		EC (DPH)	0.60	0.60	0.60
		FC (PH)	1.00	1.00	1.00
kernel		EC (DPH)	0.75	0.90	0.60
		FC (PH)	1.00	1.00	1.00
AUC		EC (DPH)	0.50	0.60	0.40
		FC (PH)	1.00	1.00	1.00
IPF		EC (DPH)	0.50	0.60	0.40
		FC (PH)	1.00	1.00	1.00

D2. Full versus masked connectivity networks

An additional analysis has been driven by the question whether classification tasks, based on masked matrices associated to (un-)directed networks (see Section 3.2), achieve different results when applied to the full matrices. We compare in Table D.7 full and masked-based classification results; the results suggest that the masking does not have a systematic detrimental effect on the analysis outcomes.

D3. 0-Dimensional Features

In this analysis carried out on both EEG and fMRI data, we focus on 0-dimensional topological features. We use the features

Carlsson Coordinates (CC), Persistent Entropy (PE) and the kernel feature of Reininghaus et al. (2015). Additionally, we investigate the 0-dimensional features Area Under the Curve (AUC) Gracia-Tabuenca et al. (2019) and Integrated Persistence Feature (IPF) Kuang et al. (2019a), as both have shown to provide good differentiations when applied to attention-deficit/hyperactivity disorder (ADHD) and Alzheimers disease, respectively. We show in Table D.8 the results on the fMRI (schizophrenia) dataset and on the scalp EEG (epilepsy) data; we present in Table D.9 the results in dimension 0 for both the scalp EEG and intracranial iEEG datasets, when using the alternative approach discussed in Section 4.2.2 and Section 4.3.

Table D.9

SVM classification with 0-dimensional topological features, for the masked FC and EC. Epil EEG and Epil iEEG denote the EEG epilepsy data for the group analysis for the scalp and intracranial data, respectively. Features (CC) and (PE) refer to the topological features of Section 3.5, (kernel), (AUC) and (IPF) are additional features discussed in the text. Reported the undirected approach (PH) or directed approach (DPH). Acc stands for Accuracy, Sns for Sensitivity and Spc for Specificity.

Data	Features	Method	Acc	Sns	Spc
Epil EEG	PE	(DPH)	0.44	0.44	0.44
		(PH)	1.00	1.00	1.00
	CC	(DPH)	0.83	0.83	0.83
Epil iEEG	PE	(PH)	0.92	0.89	0.94
		(DPH)	0.62	0.62	0.62
	CC	(PH)	0.94	0.94	0.94
	CC	(DPH)	0.59	0.56	0.62
		(PH)	0.81	0.81	0.81

Table E.10

Epilepsy: EEG analysis with variation, full report. SVM classification with 0-dimensional and 1-dimensional topological features, for the full/masked FC and full/masked EC derived features associated to the EEG dataset, following the variation of the pipeline. The features (CC) and (PE) refer to the topological features of Section 3.5, (kernel), (AUC) and (IPF) are additional features discussed in the appendix.

Feature	Full-FC	Masked-FC	Full-EC	Masked-EC	
DimAUC	0.94	1.00	0.61	0.67	
0	IPF	0.89	1.00	0.67	
	PE	0.89	1.00	0.61	
	CC	0.94	0.92	0.31	0.83
	kernel	0.92	0.92	0.50	0.53
DimPI	0.86	0.83	0.58	0.86	
1	PE	0.33	0.06	0.50	0.89
	PL	0.83	0.75	0.56	0.58
	CC	0.86	0.92	0.75	0.61
	kernel	0.83	0.81	0.69	0.86

Appendix E. Variation on EEG and iEEG dataset: the complete report

We report here the full analysis based on the variation of the pipeline (as described in Section 3.6) in the case of the EEG/iEEG recordings. In the results of Section 4.2.2 and Section 4.3, we have only shown the classification results based on the masked FC and on the masked EC matrices (see Section 3.2). For the sake of completeness, we show below in Tables E.10 and E.11 the results achieved by the topological features PL, PI, PE, CC (see Section 3.5), kernel (introduced in Appendix D.1),

and by the features AUC and IPF (as in the previous section), in both dimension 0 and dimension 1 and for Full-FC, Masked-FC, Full-EC and Masked EC.

Appendix F. Comparison of PH and DPH classification

A direct comparison between the PH and DPH performances is not straightforward given the number and intricate dependence of the classification tasks. In Fig. F.9, we present a scatter-plot where each point represents a classification of one of the datasets using one particular TDA feature (vector). The X-coordinate shows the classification accuracy obtained with DPH (computed on EC), and the Y-coordinate shows

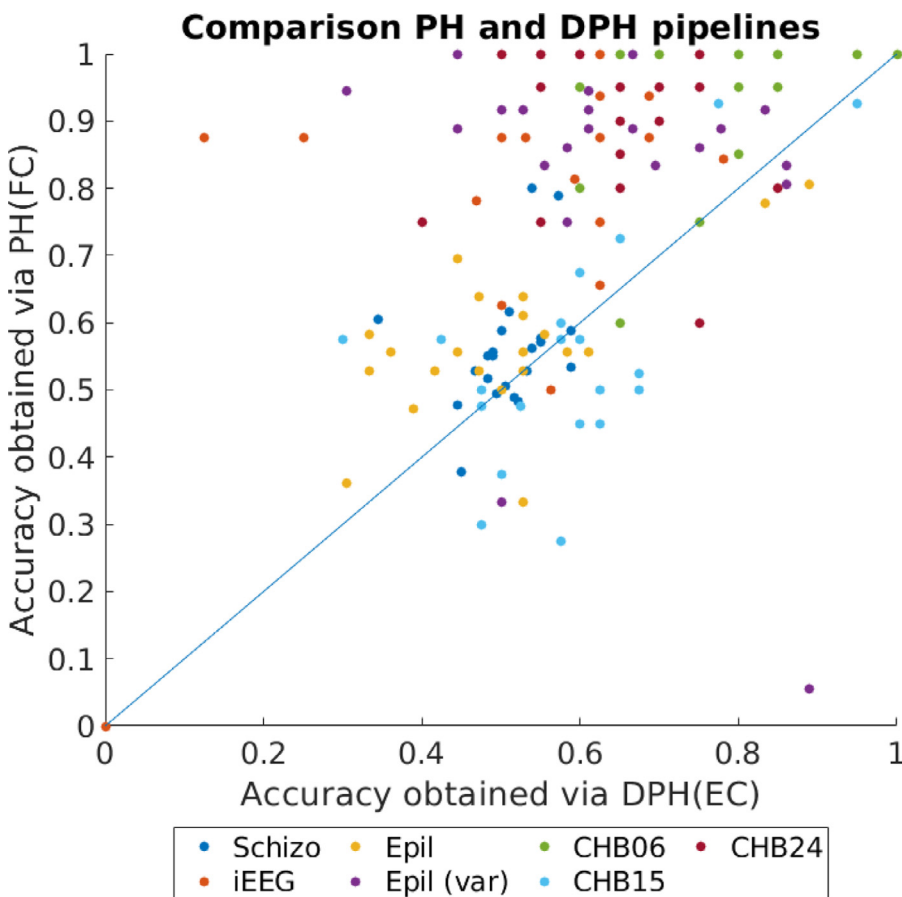


Fig. F.9. Scatter plot of classification accuracy reached using PH and DPH pipelines.

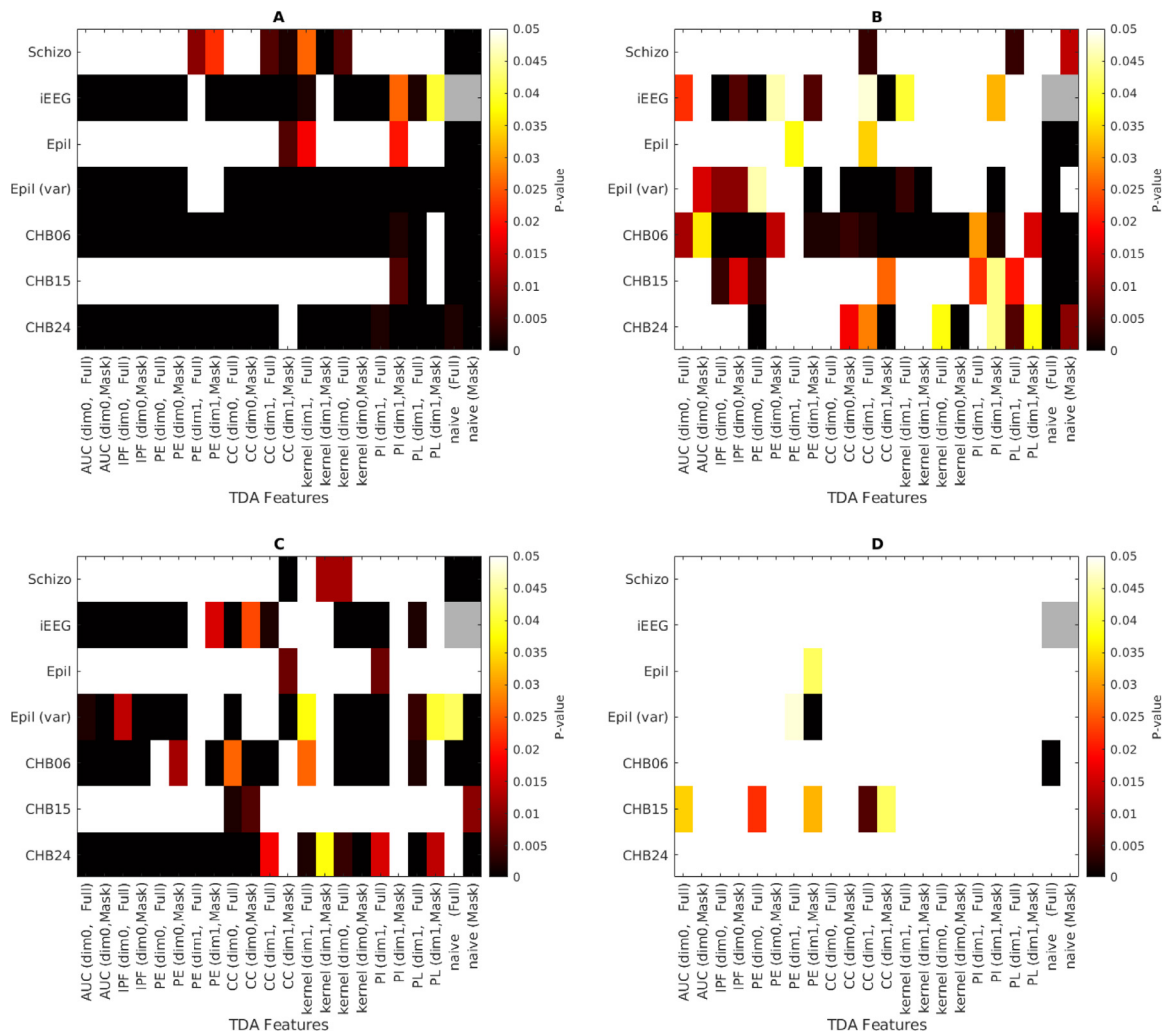


Fig. F.10. A. Uncorrected p -values for significance of classifications using FC(PH) pipeline. Null hypothesis: FC(PH)-based accuracy is random. Alternative: FC(PH)-based accuracy is higher than random. B. Uncorrected p -values for significance of classifications using EC(DPH) pipeline. Null hypothesis: EC(DPH)-based accuracy is random. Alternative: EC(DPH)-based accuracy is higher than random. C. Uncorrected p -values for comparison of classifications using FC(PH) and EC(DPH) pipeline. Null hypothesis: accuracy of classifiers are the same. Alternative hypothesis: FC(PH)-based accuracy is higher. D. Uncorrected p -values for comparison of classifications using FC(PH) and EC(DPH) pipeline. Null hypothesis: accuracy of classifiers are the same. Alternative hypothesis: EC(DPH)-based accuracy is higher. Naive classification for the iEEG dataset is not possible due to different matrix sizes across subjects.

Table E.11

Epilepsy: iEEG analysis, full report, following the variation of the pipeline. SVM classification with 0-dimensional and 1-dimensional topological features, for the full/masked FC and full/masked EC derived features. The features (CC) and (PE) refer to the topological features of Section 3.5, (kernel), (AUC) and (IPF) are additional features discussed in the appendix.

Feature	Full-FC	Masked-FC	Full-EC	Masked-EC
DimAUC	1.00	0.87	0.62	0.12
0 IPF	1.00	0.94	0.75	0.69
PE	1.00	0.94	0.75	0.62
CC	0.87	0.81	0.25	0.59
kernel	0.87	0.87	0.53	0.50
DimPI	0.81	0.66	0.59	0.62
1 PE	0.37	0.87	0.50	0.69
PL	0.78	0.62	0.47	0.50
CC	0.87	0.84	0.62	0.78
kernel	0.75	0.50	0.62	0.56

the corresponding PH (computed on FC) results. Different datasets are marked with different colours.

We observe that, in some cases, PH and DPH are comparable - shown by the points laying close to the diagonal; in other cases, PH seems to work better - upper triangle points. The cloud around (0.5,0.5) corresponds to results close to random for both input matrices, mainly for the schizophrenia and group epilepsy studies.

To statistically evaluate the two approaches, we present in Fig. F.10 the results of permutation tests. Firstly, we assess the classification significance of both pipelines, PH- and DPH-based, by permuting the obtained labels of data points. By repeating this 500 times (500 permutations) we get a distribution of random classifications; we use accuracy as the test statistic, to compute the p -value for the null hypothesis that the accuracy of our classifier is not distinguishable from random. Corresponding (uncorrected) p -values can be found in Fig. F.10, top sub-panel.

For the comparison between the DPH- and PH-based classifications, we performed another permutation test. Now, for every tested dataset and TDA feature, we work with pairs of vectors: labels for the classifier that we get from PH and from DPH. For every data point, we either flip the corresponding coordinate between these vectors, or leave it with equal probability. In this way, we generate surrogates representing the assumption “there is no difference between the PH and DPH approach

performance". We test the mentioned null hypothesis against two alternatives: PH is better than DPH, and DPH is better than PH. For that we use test statistic given by the accuracy of PH minus the accuracy of DPH, concentrating on its high positive and negative values for testing the mentioned two alternative hypotheses, respectively. Corresponding (uncorrected) p -values can be found in Fig. F.10, bottom subpanel.

From the Fig. F.10 it seems that the most important parameter, influencing classification, is the dataset. For our datasets, schizophrenia, epilepsy on a group level and patient CHB15 are those where the TDA accuracies are most random (not the same for the naive approaches). In the other cases, both PH and DPH seem to work better than random. We present uncorrected p -values since due to the complexity of dependencies it is difficult to choose an appropriate multiple testing correction. Both the TDA features and datasets are not independent, so it is difficult to compute the degrees of freedom for the correction. However, the number of significant classifications for both PH and DPH is clearly substantially higher than the expected 5% of false positives; so, generally, we can conclude that both pipelines work in favorable conditions. Furthermore, testing the "PH is better than DPH" alternative, seems to provide a larger number of significant classifications, which would suggest that generally PH performs better than DPH.

References

Adams, H., Emerson, T., Kirby, M., Neville, R., Peterson, C., Shipman, P., Chepushanova, S., Hanson, E., Motta, F., Ziegelmeier, L., 2017. Persistence images: a stable vector representation of persistent homology. *J. Mach. Learn. Res.* 18 (8), 1–35.

Adcock, A., Carlsson, E., Carlsson, G., 2013. The ring of algebraic functions on persistence bar codes. *Homol. Homotopy Appl.* 18. doi:10.4310/HHA.2016.v18.n1.a21.

Aktas, M.E., Akbas, E., Fatmaoui, A.E., 2019. Persistence homology of networks: methods and applications. *Appl. Netw. Sci.* 4 (61). doi:10.1007/s41109-019-0179-3.

Barnett, L., Seth, A., 2013. The MVGC multivariate Granger causality toolbox: a new approach to Granger-causal inference. *J. Neurosci. Methods* 223. doi:10.1016/j.jneumeth.2013.10.018.

Bendich, P., Marron, J., Miller, E., Pieloch, A., Skwerer, S., 2014. Persistent homology analysis of brain artery trees. *Ann. Appl. Stat.* 10. doi:10.1214/15-AOAS886.

Brodzki, J., Belchí, F., Djukanovic, R., Conway, J., Pirashvili, M., Bennett, M., 2018. Lung topology characteristics in patients with chronic obstructive pulmonary disease. *Sci. Rep.* 8. doi:10.1038/s41598-018-23424-0.

Bubenik, P., Dlotko, P., 2017. A persistence landscapes toolbox for topological statistics. *J. Symb. Comput.* 78, 91–114. doi:10.1016/j.jsc.2016.03.009. Algorithms and Software for Computational Topology

Bullmore, E., Sporns, O., 2009. Complex brain networks: graph theoretical analysis of structural and functional systems. *Nat. Rev. Neurosci.* 10 (3), 186–198. doi:10.1038/nrn2575.

Carlsson, G., 2009. Topology and data. *Bull. Amer. Math. Soc.* 46, 255–308. doi:10.1090/S0273-0979-09-01249-X.

Castro, E., Martínez-Ramón, M., Pearlson, G., Sui, J., Calhoun, V.D., 2011. Characterization of groups using composite kernels and multi-source fMRI analysis data: application to schizophrenia. *NeuroImage* 58 (2), 526–536.

Chazal, F., Michel, B., 2017. An introduction to topological data analysis: fundamental and practical aspects for data scientists. *arXiv preprint arXiv:1710.04019*.

Choi, H., Kim, Y., Kang, H., Hyekyoung, L., Im, H.-J., Hwang, D., Kim, E., Chung, J.-K., Lee, D., 2014. Abnormal metabolic connectivity in the pilocarpine-induced epilepsy rat model: a multiscale network analysis based on persistent homology. *NeuroImage* 99. doi:10.1016/j.neuroimage.2014.05.039.

Chowdhury, S., Mémoli, F., 2018. Persistent path homology of directed networks. pp. 1152–1169. 10.1137/1.9781611975031.75

Chung, M.K., Hanson, J.L., Ye, J., Davidson, R.J., Pollak, S.D., 2014. Persistent homology in sparse regression and its application to brain morphometry. *IEEE Trans. Med. Imaging* 34. doi:10.1109/TMI.2015.2416271.

Cohen-Steiner, D., Edelsbrunner, H., Harer, J., 2007. Stability of persistence diagrams. *Discrete Comput. Geom.* 37, 103–120.

Dey, T.K., Mandal, S., Varcho, W., 2017. Improved image classification using topological persistence. In: Hullin, M., Klein, R., Schultz, T., Yao, A. (Eds.), *Vision, Modeling and Visualization*. The Eurographics Association doi:10.2312/vmv.20171272.

Dunaeva, O., Edelsbrunner, H., Lukyanov, A., Machin, M., Malkova, D., Kuvaev, R., Kashin, S., 2016. The classification of endoscopic images with persistent homology. *Pattern Recognit. Lett.* 83, 13–22. doi:10.1016/j.patrec.2015.12.012. Geometric, topological and harmonic trends to image processing

Edelsbrunner, H., Letscher, D., Zomorodian, A., 2002. Topological persistence and simplification. *Discrete Comput. Geom.* 28, 511–533.

Ellis, C.T., Lesnick, M., Henselman-Petrusek, G., Keller, B., Cohen, J.D., 2019. Feasibility of topological data analysis for event-related fMRI. *Netw. Neurosci.* 3 (3), 695–706. doi:10.1162/netn_a_00095.

Friston, K.J., 1994. Functional and effective connectivity in neuroimaging: a synthesis. *Hum. Brain Mapp.* 2 (1–2), 56–78. doi:10.1002/hbm.460020107.

Friston, K.J., 2011. Functional and effective connectivity: a review. *Brain Connect.* 1 (1), 13–36. doi:10.1089/brain.2011.0008. PMID: 22432952

Frosini, P., 1990. A distance for similarity classes of submanifolds of a Euclidean space. *Bull. Aust. Math. Soc.* 42 (3), 407–415. doi:10.1017/S0004972700028574.

Gallos, L.K., Makse, H.A., Sigman, M., 2012. A small world of weak ties provides optimal global integration of self-similar modules in functional brain networks. *Proc. Natl. Acad. Sci.* 109 (8), 2825–2830. doi:10.1073/pnas.1106612109. <https://www.pnas.org/content/109/8/2825.full.pdf>

Garg, A., Lu, D., Popuri, K., Beg, M.F., 2016. Cortical geometry network and topology markers for Parkinson's disease. *ArXiv*.

Giansiracusa, N., Giansiracusa, R., Moon, C., 2019. Persistent homology machine learning for fingerprint classification. pp. 1219–1226. 10.1109/ICMLA.2019.00201

Gidea, M., 2017. Topology data analysis of critical transitions in financial networks. *SSRN Electron. J.* doi:10.2139/ssrn.2903278.

Gracia-Tabuenca, Z., Díaz-Patiño, J.C., Arelio, I., Alcauter, S., 2019. Topological data analysis reveals robust alterations in the whole-brain and frontal lobe functional connectomes in attention-deficit/hyperactivity disorder. *bioRxiv* doi:10.1101/751008.

Hartman, D., Hlinka, J., Palu, M., Mantini, D., Corbetta, M., 2011. The role of nonlinearity in computing graph-theoretical properties of resting-state functional magnetic resonance imaging brain networks. *Chaos* 21 (1), 013119. doi:10.1063/1.3553181.

Hatcher, A., 2000. *Algebraic Topology*. Cambridge Univ. Press, Cambridge.

Hlinka, J., Hartman, D., Vejmelka, M., Runge, J., Marwan, N., Kurths, J., Paluš, M., 2013. Reliability of inference of directed climate networks using conditional mutual information. *Entropy* 15 (6), 2023–2045. doi:10.3390/e15062023.

Hlinka, J., Palus, M., Vejmelka, M., Mantini, D., Corbetta, M., 2011. Functional connectivity in resting-state fMRI: is linear correlation sufficient? *NeuroImage* 54, 2218–2225. doi:10.1016/j.neuroimage.2010.08.042.

Horak, D., Maletić, S., Rajković, M., 2009. Persistent homology of complex networks. *J. Stat. Mech.* (03).

Khalid, A., Kim, B.S., Chung, M.K., Ye, J.C., Jeon, D., 2014. Tracing the evolution of multi-scale functional networks in a mouse model of depression using persistent brain network homology. *NeuroImage* 101. doi:10.1016/j.neuroimage.2014.07.040.

Kopal, J., Pidnebesna, A., Tomeček, D., Tintěra, J., Hlinka, J., 2020. Typicality of functional connectivity robustly captures motion artifacts in rs-fMRI across datasets, atlases and preprocessing pipelines. *Hum. Brain Mapp.* Submitted

Kořenek, J., Hlinka, J., 2020. Causal network discovery by iterative conditioning: comparison of algorithms. *Chaos* 30 (1), 013117. doi:10.1063/1.5115267.

Kuang, L., Han, X., Chen, H.C., Reiman, R.J., Eric, M., W.Y., 2019. A concise and persistent feature to study brain resting-state network dynamics: findings from the Alzheimer's disease neuroimaging initiative. *Hum. Brain Mapp.* 40 (4), 1062–1081. doi:10.1002/hbm.24383.

Kuang, L., Zhao, D., Xing, J., Chen, Z., Xiong, F., Han, X., 2019. Metabolic brain network analysis of FDG-PET in Alzheimers disease using kernel-based persistent features. *Molecules* 24, 2301. doi:10.3390/molecules24122301.

Lee, D., 2019. Clinical personal connectomics using hybrid PET/MRI. *Nucl. Med. Mol. Imaging* 53, 1–11. doi:10.1007/s13139-019-00572-3.

Lee, H., Chung, M.K., Kang, H., Kim, B.-N., Lee, D.S., 2011. Computing the shape of brain networks using graph filtration and Gromov-Hausdorff metric. In: Fichtinger, G., Martel, A., Peters, T. (Eds.), *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2011*. Springer Berlin Heidelberg, Berlin, Heidelberg, pp. 302–309.

Lee, H., Chung, M.K., Kang, H., Kim, B.-N., Lee, D.S., 2011. Discriminative persistent homology of brain networks. In: 2011 International Symposium on Biomedical Imaging, pp. 841–844. doi:10.1109/ISBI.2011.5872535.

Lee, H., Kang, H., Chung, M.K., Kim, B., Lee, D.S., 2012. Persistent brain network homology from the perspective of dendrogram. *IEEE Trans. Med. Imaging* 31 (12), 2267–2277.

Liao, W., Zhang, Z., Pan, Z., Mantini, D., Ding, J., Duan, X., Luo, C., Lu, G., Chen, H., 2010. Altered functional connectivity and small-world in mesial temporal lobe epilepsy. *PLoS ONE* 5 (1).

Lutkepohl, H., 2007. *New Introduction to Multiple Time Series Analysis*. Springer Publishing Company, Incorporated.

Lynall, M.-E., Bassett, D.S., Kerwin, R., McKenna, P.J., Kitzbichler, M., Muller, U., Bullmore, E., 2010. Functional connectivity and brain networks in schizophrenia. *J. Neurosci.* 30 (28), 9477–9487. doi:10.1523/JNEUROSCI.0333-10.2010.

Lütgehetmann, D., Govc, D., Smith, J.P., Levi, R., 2020. Computing persistent homology of directed flag complexes. *Algorithms* 13 (1), 19. doi:10.3390/a13010019.

Masulli, P., Villa, A., 2016. The topology of the directed clique complex as a network invariant. *SpringerPlus* 5.

Merelli, E., Piangerelli, M., Rucco, M., Toller, D., 2016. A topological approach for multivariate time series characterization: the epileptic brain. *EAI Endorsed Trans. Self-Adaptive Syst.* 2 (7). doi:10.4108/eai.3-12-2015.2262525.

Mikolas, P., Hlinka, J., Skoch, A., Pitra, Z., Frodl, T., Spaniel, F., Hajek, T., 2018. Machine learning classification of first-episode schizophrenia spectrum disorders and controls using whole brain white matter fractional anisotropy. *BMC Psychiatry* 18 (1).

Munkres, J.R., 1984. *Elements of Algebraic Topology*. Addison Wesley Publishing Company, Inc., 2725 Sand Hill Road Menlo Park, California 94025.

Oliver, I., Hlinka, J., Kopal, J., Runge, J., Daviden, J., 2019. Quantifying the variability in resting-state networks. *Entropy* 21 (9), 882. doi:10.3390/e21090882.

Otter, N., Porter, M.A., Tillmann, U., Grindrod, P., Harrington, H.A., 2017. A roadmap for the computation of persistent homology. *EPJ Data Sci.* 6 (17).

Pachauri, D., Hinrichs, C., Chung, M.K., Johnson, S.C., Singh, V., 2011. Topology-based kernels with application to inference problems in Alzheimer's disease. *IEEE Trans. Med. Imaging* 30, 1760–1770. doi:10.1109/TMI.2011.2147327.

Park, H.-J., Friston, K., 2013. Structural and functional brain networks: from connections to cognition. *Science* 342 (6158), 1238411. doi:10.1126/science.1238411.

Petri, G., Sciamiero, M., Donato, I., Vaccarino, F., 2013. Topological strata of weighted complex networks. *PLoS ONE* 8 (6), 1–8. doi:10.1371/journal.pone.0066506.

- Phinyomark, A., Ibáñez-Marcelo, E., Petri, G., 2017. Resting-state fMRI functional connectivity: big data preprocessing pipelines and topological data analysis. *IEEE Trans. Big Data* 3 (4), 415–428. doi:[10.1109/TBDATA.2017.2734883](https://doi.org/10.1109/TBDATA.2017.2734883).
- Piangerelli, M., Rucco, M., Tesei, L., Merelli, E., 2018. Topological classifier for detecting the emergence of epileptic seizures. *BMC Res. Notes* 11. doi:[10.1186/s13104-018-3482-7](https://doi.org/10.1186/s13104-018-3482-7).
- Rathore, A., Palande, S., Anderson, J., Zielinski, B., Fletcher, P., Wang, B., 2019. Autism classification using topological features and deep learning: a cautionary tale. pp. 736–744. [10.1007/978-3-030-32248-9_82](https://doi.org/10.1007/978-3-030-32248-9_82)
- Reimann, M.W., Nolte, M., Scolamiero, M., Turner, K., Perin, R., Chindemi, G., Dłotko, P., Levi, R., Hess, K., Markram, H., 2017. Cliques of neurons bound into cavities provide a missing link between structure and function. *Front. Comput. Neurosci.* 11, 48. doi:[10.3389/fncom.2017.00048](https://doi.org/10.3389/fncom.2017.00048).
- Reininghaus, J., Huber, S., Bauer, U., Kwitt, R., 2015. A stable multi-scale kernel for topological machine learning. Boston, MA, pp. 4741–4748. [10.1109/CVPR.2015.7299106](https://doi.org/10.1109/CVPR.2015.7299106)
- Rubin-Falcone, H., Zanderigo, F., Thapa-Chhetry, B., Lan, M., Miller, J.M., Sublette, M.E., Oquendo, M.A., Hellerstein, D.J., McGrath, P.J., Stewart, J.W., Mann, J.J., 2018. Pattern recognition of magnetic resonance imaging-based gray matter volume measurements classifies bipolar disorder and major depressive disorder. *J. Affect. Disord.* 227, 498–505.
- Rucco, M., Castiglione, F., Merelli, E., Pettini, M., 2014. Characterisation of the idiotypic immune network through persistent entropy. In: *Proceedings of ECCS 2014*, Lucca, Italy, pp. 117–128. doi:[10.1007/978-3-319-29228-1_11](https://doi.org/10.1007/978-3-319-29228-1_11).
- Siddiqui, M.K., Morales-Menendez, R., Huang, X., Hussain, N., 2020. A review of epileptic seizure detection using machine learning classifiers. *Brain Inf.* 7 (1), 5. doi:[10.1186/s40708-020-00105-1](https://doi.org/10.1186/s40708-020-00105-1).
- Smith, S.M., 2012. The future of fMRI connectivity. *NeuroImage* 62 (2), 1257–1266. doi:[10.1016/j.neuroimage.2012.01.022](https://doi.org/10.1016/j.neuroimage.2012.01.022). 20 YEARS OF fMRI
- Solo, V., Poline, J., Lindquist, M.A., Simpson, S.L., Bowman, F.D., Chung, M.K., Cassidy, B., 2018. Connectivity in fMRI: blind spots and breakthroughs. *IEEE Trans. Med. Imaging* 37 (7), 1537–1550.
- Stolz, B. J., Emerson, T., Nahkuri, S., Porter, M. A., Harrington, H. A., 2018. Topological data analysis of task-based fMRI data from experiments on schizophrenia. *arXiv preprint arXiv:1809.08504*.
- Tauzin, G., Lupo, U., Tunstall, L., Pérez, J. B., Caorsi, M., Medina-Mardones, A., Dassatti, A., Hess, K., 2020. giotto-tda: a topological data analysis toolkit for machine learning and data exploration. *arXiv preprint arXiv:2004.02551*.
- Telesford, Q., Simpson, S., Burdette, J., Hayasaka, S., Laurienti, P., 2011. The brain as a complex system: using network science as a tool for understanding the brain. *Brain Connect.* 1, 295–308. doi:[10.1089/brain.2011.0055](https://doi.org/10.1089/brain.2011.0055).
- Tralie, C., Saul, N., Bar-On, R., 2018. Ripser.py: a lean persistent homology library for Python. *J. Open Source Softw.* 3 (29), 925. doi:[10.21105/joss.00925](https://doi.org/10.21105/joss.00925).
- Turner, K., 2019. Rips filtrations for quasimetric spaces and asymmetric functions with stability results. *Algebr. Geom. Topol.* 19 (3), 1135–1170. doi:[10.2140/agt.2019.19.1135](https://doi.org/10.2140/agt.2019.19.1135).
- Valdes-Sosa, P.A., Roebroeck, A., Daunizeau, J., Friston, K., 2011. Effective connectivity: influence, causality and biophysical modeling. *NeuroImage* 58 (2), 339–361. doi:[10.1016/j.neuroimage.2011.03.058](https://doi.org/10.1016/j.neuroimage.2011.03.058). <http://www.sciencedirect.com/science/article/pii/S1053811911003375>
- Wang, Y., Ombao, H., Chung, M., 2015. Topological seizure origin detection in electroencephalographic signals. In: *2015 IEEE International Symposium on Biomedical Imaging*, pp. 351–354. doi:[10.1109/ISBI.2015.7163885](https://doi.org/10.1109/ISBI.2015.7163885).
- Wasserman, L., 2018. Topological data analysis. *Annu. Rev. Stat. Appl.* 5 (1), 501–532. doi:[10.1146/annurev-statistics-031017-100045](https://doi.org/10.1146/annurev-statistics-031017-100045).
- Wong, E., Palande, S., Wang, B., Zielinski, B., Anderson, J., Fletcher, P., 2016. Kernel partial least squares regression for relating functional brain network topology to clinical measures of behavior. In: *IEEE 13th International Symposium on Biomedical Imaging (ISBI)*, pp. 1303–1306. doi:[10.1109/ISBI.2016.7493506](https://doi.org/10.1109/ISBI.2016.7493506).
- Yesilli, M. C., Khasawneh, F. A., Otto, A., 2019. Topological feature vectors for chatter detection in turning processes. *arXiv preprint arXiv:1905.08671*.
- Zomorodian, A., Carlsson, G., 2005. Computing persistent homology. *Discrete Comput. Geom.* 33, 249–274.
- Zurita, M., Montalba, C., Labbé, T., Cruz, J.P., Dalboni da Rocha, J., Tejos, C., Ciampi, E., Cárcamo, C., Sitaram, R., Uribe, S., et al., 2018. Characterization of relapsing-remitting multiple sclerosis patients using support vector machine classifications of functional and diffusion MRI data. *NeuroImage Clin.* 20, 724–730. doi:[10.1016/j.nicl.2018.09.002](https://doi.org/10.1016/j.nicl.2018.09.002).