HELMHOLTZ AIH Institute of AI for Health

### **Topological Machine Learning: The (W) Hole Truth Lecture 5** Bastian Rieck (@Pseudomanifold)

### **Preliminaries**

Do you have feedback or any questions? Write to bastian.rieck@helmholtz-muenchen.de or reach out to @Pseudomanifold on Twitter. You can find the slides and additional information with links to more literature here:

https://heidelberg.topology.rocks

- the persistence diagram is the 'basic' topological feature descriptor.
- ☆ Multiple alternatives exist, with different key properties.
- ☆ Their choice is application-dependent.

### In this lecture

Putting everything together



#### How can we build topology-based machine learning models?

Topological machine learning



Point cloud



Point cloud

Persistent homology



Point cloud

Persistent homology

Persistence diagram(s)





- A. Poulenard, P. Skraba and M. Ovsjanikov, 'Topological Function Optimization for Continuous Shape Matching', *Computer Graphics Forum* 37.5, 2018, pp. 13–25
- M. Moor\*, M. Horn\*, B. Rieck<sup>†</sup> and K. Borgwardt<sup>†</sup>, 'Topological Autoencoders', Proceedings of the 37th International Conference on Machine Learning (ICML), 2020, pp. 7045–7054, arXiv: 1906.00722 [cs.LG]
- A. Carrière, F. Chazal, M. Glisse, Y. Ike, H. Kannan and Y. Umeda, 'Optimizing persistent homology based functions', *Proceedings of the 38th International Conference on Machine Learning (ICML)*, 2021, pp. 1294–1303

Finding topology-based representations

1. We detering a new topological loss term for anime

#### Tepological Autoencoder

#### ichael Moor 117 May Bern 117 Bastian Birch 117 Karsien Berrmardt

ing a special simplicial complex, the Veteria Equ



Michael Moor ♥ Michael\_D\_Moor



Max Horn ♥ ExpectationMax



Karsten Borgwardt ♥ kmborgwardt

M. Moor<sup>\*</sup>, M. Horn<sup>\*</sup>, **B. Rieck**<sup> $\dagger$ </sup> and K. Borgwardt<sup> $\dagger$ </sup>, 'Topological Autoencoders', *Proceedings of the 37th International* Conference on Machine Learning (ICML), 2020, pp. 7045–7054, arXiv: 1906.00722 [cs.LG]

Motivation

Motivation, continued

Overview



Main intuition

Align persistence diagrams of an *input batch* and of a *latent batch* using a loss function!

### Why this works in theory

Let X be a point cloud of cardinality n and  $X^{(m)}$  be one subsample of X of cardinality m, i.e.  $X^{(m)} \subseteq X$ , sampled without replacement. We can bound the probability of the persistence diagrams of  $X^{(m)}$  exceeding a threshold in terms of the bottleneck distance as

$$\mathbb{P}\Big(\mathsf{W}_{\!\infty}\!\left(\mathcal{D}^{X}\!,\mathcal{D}^{X^{(m)}}\right)\!>\!\epsilon\Big)\leq\mathbb{P}\Big(\mathsf{dist}_{\mathsf{H}}\!\left(X,X^{(m)}\right)\!>\!2\epsilon\Big),$$

where dist<sub>H</sub> denotes the Hausdorff distance. In other words: *mini-batches are topologically similar if the subsampling is not too coarse.* 

Gradient calculation intuition

Distance matrix **A**  $\begin{bmatrix} 0 & 1 & 9 & 10 \\ 1 & 0 & 7 & 8 \\ 9 & 7 & 0 & 3 \\ 10 & 8 & 3 & 0 \end{bmatrix}$ 

Every point in the persistence diagram can be mapped to *one* entry in the distance matrix! Each entry *is* a distance, so it can be changed during training (at least in the latent space).

Gradient calculation intuition



Every point in the persistence diagram can be mapped to *one* entry in the distance matrix! Each entry *is* a distance, so it can be changed during training (at least in the latent space).

Gradient calculation intuition



Every point in the persistence diagram can be mapped to *one* entry in the distance matrix! Each entry *is* a distance, so it can be changed during training (at least in the latent space).

Loss term

$$\mathcal{L}_t := \mathcal{L}_{\mathcal{X} \to \mathcal{Z}} + \mathcal{L}_{\mathcal{Z} \to \mathcal{X}}$$

 $\mathcal{L}_{\mathcal{X}\to\mathcal{Z}} := \frac{1}{2} \left\| \mathbf{A}^X \left[ \pi^X \right] - \mathbf{A}^Z \left[ \pi^X \right] \right\|^2 \qquad \qquad \mathcal{L}_{\mathcal{Z}\to\mathcal{X}} := \frac{1}{2} \left\| \mathbf{A}^Z \left[ \pi^Z \right] - \mathbf{A}^X \left[ \pi^Z \right] \right\|^2$ 

- 🕸 🛛  $\mathcal{Z}$ : latent space
- $\diamond$  **A**<sup>X</sup>: distances in input mini-batch
- $\Rightarrow$  **A**<sup>Z</sup>: distances in latent mini-batch
- $\hat{\pi}^Z$ : persistence pairing of latent mini-batch

The loss is bi-directional!

## **Qualitative evaluation**

'Spheres' data set



### **Quantitative evaluation**

Data set	Method	$KL_{0.01}$	$KL_{0.1}$	$KL_1$	$\ell$ -MRRE	$\ell ext{-Cont}$	$\ell$ -Trust	$\ell$ -RMSE	MSE (data)
	Isomap	0.181	0.420	0.00881	0.246	0.790	0.676	10.4	
	PCA	0.332	0.651	0.01530	0.294	0.747	0.626	11.8	0.9610
'Sphoros'	t-SNE	0.152	0.527	0.01271	0.217	0.773	0.679	<u>8.1</u>	
Spheres	UMAP	0.157	0.613	0.01658	0.250	0.752	0.635	9.3	
	AE	0.566	0.746	0.01664	0.349	0.607	0.588	13.3	0.8155
	ТороАЕ	0.085	0.326	0.00694	0.272	0.822	0.658	13.5	0.8681
	PCA	0.356	0.052	0.00069	0.057	0.968	0.917	9.1	0.1844
	t-SNE	0.405	0.071	0.00198	0.020	0.967	0.974	41.3	
'Fashion-MNIST'	UMAP	0.424	0.065	0.00163	0.029	0.981	0.959	13.7	
	AE	0.478	0.068	0.00125	0.026	0.968	0.974	20.7	0.1020
	ТороАЕ	0.392	0.054	0.00100	0.032	0.980	0.956	20.5	0.1207
'MNIST'	PCA	0.389	0.163	0.00160	0.166	0.901	0.745	13.2	0.2227
	t-SNE	0.277	0.133	0.00214	0.040	0.921	0.946	22.9	
	UMAP	0.321	0.146	0.00234	0.051	0.940	0.938	14.6	
	AE	0.620	0.155	0.00156	0.058	0.913	0.937	18.2	0.1373
	ТороАЕ	0.341	<u>0.110</u>	0.00114	0.056	0.932	0.928	19.6	0.1388

# Topology-driven graph learning

Using 'classical' machine learning models

- Calculate degree filtration (or another descriptor)
- 2 Repeat the analysis pipeline described above
- 3 Learn weights for topological descriptors to improve predictive power<sup>1</sup>

<sup>1</sup>Q. Zhao and Y. Wang, 'Learning metrics for persistence-based summaries and applications for graph classification', *Advances in Neural Information Processing Systems*, ed. by H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox and R. Garnett, vol. 32, Curran Associates, Inc., 2019, pp. 9855–9866

HELMHOLTZ MUNICI<del>)</del>

### **Betti curves**

Classification scenario example

- Use REDDIT-BINARY data set (co-occurrence graphs)
- ☆ Calculate filtration based on vertex degree
- pprox Calculate persistence diagrams for d=1 (cycles)











### Static topological features for graph classification

A Persistent Weideiler-Lehman Procedure for Graph Classification

Basian Birch<sup>11</sup> Christian Boch<sup>11</sup> Karsten Borgwordt<sup>1</sup>

#### Abrill

The Workshop Learns is hole poply density comparising spheroscenes is hole poply density contain takes. However, its solution faintees and pople, topological batters lawns for damages and the se approxcess of the second spheroscenes and pople, topological data and the second spheroscenes from averaging apple into earlier and teams from the second spheroscenes with handling of the second spheroscenes and topological distances in default contents with handling of the second spheroscenes and pople, the second spheroscenes are applied with the second spheroscenes and administical formation is default for the second population of the spheroscenes and administration of the shift of the schemation as any general administical of the shift of the schemation areas and administical the schemation as any second administical the schemation and schematical the schemation and the schematical the schemation and schematical the schematical the schemation administical the schemation and schematical the schematical the schemation and schematical schematical the schemat

1. Introduction

ménerel data sen are abliquiteux in a vasiety application demains, each of them proing a su dauge while also requiring different index to k comment talk introduce graph classification of assists of methods action. These methods

version of ed. 2010), we observe that an analysis (Manuch 2017) to all armitichations are as follows: (b) the later address for a measure file reveause of ed. in anomare more version with the second set of PUTS, which makes the restrict any provide a set of the all second set of PUTS.

 annu a series of it with the series for any first test we show in the a presendence of version of the angles of easy.
 We develop a topology hand kernel that new antimation variant of the WL subditiation procedure to classify nonativitized pople.
 We demonstrate that one proposed features perform

ad contribution. "Department of Responses, Keisener and inclusion of Responses, Keisener and inclusion of the series. The Series Series

Proceedings of the 20<sup>40</sup> International Conference on Machine Journing, Long Banch, California, PMLB 97, 2019. Copyrigh 2019 In the authoritic



Christian Bock



Karsten Borgwardt ♥ kmborgwardt

- The Weisfeiler–Lehman algorithm vectorises labelled graphs
- Persistent homology captures relevant topological features
- ☆ We can *combine* them to obtain a *generalised* formulation
- this requires a distance between multisets

**B. Rieck**<sup>\*</sup>, C. Bock<sup>\*</sup> and K. Borgwardt, 'A Persistent Weisfeiler–Lehman Procedure for Graph Classification', *Proceedings of the 36th International Conference on Machine Learning (ICML)*, ed. by K. Chaudhuri and R. Salakhutdinov, Proceedings of Machine Learning Research 97, PMLR, 2019, pp. 5448–5458

### A distance between label multisets

Let  $A = \{l_1^{a_1}, l_2^{a_2}, \dots\}$  and  $B = \{l_1^{b_1}, l_2^{b_2}, \dots\}$  be two multisets that are defined over the same label alphabet  $\Sigma = \{l_1, l_2, \dots\}$ .

Transform the sets into count vectors, i.e.  $\vec{x} := [a_1, a_2, \dots]$  and  $\vec{y} := [b_1, b_2, \dots]$ .

Calculate their multiset distance as

$$\operatorname{dist}(\vec{x}, \vec{y}) := \left(\sum_{i} |a_i - b_i|^p\right)^{\frac{1}{p}},$$

i.e. the  $p^{th}$  Minkowski distance, for  $p \in \mathbb{R}$ . Since nodes and their multisets are in one-to-one correspondence, we now have a metric on the graph!

### **Multiset distance**

Example for p=1



$$dist(C, E) = dist(\{\bullet^3, \bullet^1\}, \{\bullet^2, \bullet^1\})$$
  
= dist([3, 1], [2, 1])  
= 1  
$$dist(C, A) = dist(\{\bullet^3, \bullet^1\}, \{\bullet^1\})$$
  
= dist([3, 1], [1, 0])  
= 3

### Extending the multiset distance to a distance between vertices

Use vertex label from *previous* Weisfeiler–Lehman iteration, i.e.  $I_{v_i}^{(h-1)}$ , as well as  $I_{v_i}^{(h)}$ , the one from the *current* iteration:

$$\operatorname{dist}(v_i, v_j) := \left[\mathsf{I}_{v_i}^{(h-1)} \neq \mathsf{I}_{v_j}^{(h-1)}\right] + \operatorname{dist}\left(\mathsf{I}_{v_i}^{(h)}, \mathsf{I}_{v_j}^{(h)}\right) + \tau$$

 $\tau \in \mathbb{R}_{>0}$  is required to make this into a proper metric. This turns *any* labelled graph into a weighted graph whose persistent homology we can calculate!









### Persistence-based Weisfeiler-Lehman feature vectors

### **Connected components**

$$\begin{split} \Phi_{\mathsf{P-WL}}^{(h)} &:= \left[ \mathfrak{p}^{(h)}(l_0), \mathfrak{p}^{(h)}(l_1), \dots \right] \\ \mathfrak{p}^{(h)}(l_i) &:= \sum_{\mathsf{I}(v) = l_i} \mathsf{pers}(v)^p, \end{split}$$

Cycles

$$\begin{split} \Phi_{\mathsf{P}\text{-WL-C}}^{(h)} &:= \left[ \mathfrak{z}^{(h)}(l_0), \mathfrak{z}^{(h)}(l_1), \ldots \right] \\ \mathfrak{z}^{(h)}(l_i) &:= \sum_{l_i \in \mathsf{I}(u,v)} \mathsf{pers}(u,v)^p, \end{split}$$

### Persistence-based Weisfeiler-Lehman feature vectors

#### **Connected components**

$$\begin{split} \Phi_{\mathsf{P-WL}}^{(h)} &:= \left[ \mathfrak{p}^{(h)}(l_0), \mathfrak{p}^{(h)}(l_1), \ldots \right] \\ \mathfrak{p}^{(h)}(l_i) &:= \sum_{\mathsf{I}(v) = l_i} \mathsf{pers}(v)^p, \end{split}$$

Cycles

$$\begin{split} \Phi_{\mathsf{P-WL-C}}^{(h)} &:= \left[\mathfrak{z}^{(h)}(l_0), \mathfrak{z}^{(h)}(l_1), \ldots\right] \\ \mathfrak{z}^{(h)}(l_i) &:= \sum_{l_i \in \mathsf{I}(u,v)} \mathsf{pers}(u,v)^p, \end{split}$$

#### Bonus

We can re-define the vertex distance to obtain the original Weisfeiler–Lehman subtree features (plus information about cycles):

$$\mathsf{dist}(v_i, v_j) := \begin{cases} 1 & \text{if } v_i \neq v_j \\ 0 & \text{otherwise} \end{cases}$$

### **Classification results**

	D & D	MUTAG	NCI1	NCI109	PROTEINS	PTC-MR	PTC-FR	PTC-MM	PTC-FM
V-Hist E-Hist	78.32 ± 0.35 72.90 ± 0.48	85.96 ± 0.27 85.69 ± 0.46	64.40 ± 0.07 63.66 ± 0.11	63.25 ± 0.12 63.27 ± 0.07	72.33 ± 0.32 72.14 ± 0.39	58.31 ± 0.27 55.82	68.13 ± 0.23 65.53	66.96 ± 0.51 61.61	57.91 ± 0.83 59.03
RetGK*	81.60 ± 0.30	90.30 ± 1.10	84.50 ± 0.20		75.80 ± 0.60	62.15 ± 1.60	67.80 ± 1.10	67.90 ± 1.40	63.90 ± 1.30
WL Deep-WL*	79.45 ± 0.38	87.26 ± 1.42 82.94 ± 2.68	85.58 ± 0.15 80.31 ± 0.46	$\begin{array}{c} 84.85 \pm 0.19 \\ 80.32 \pm 0.33 \end{array}$	76.11 ± 0.64 75.68 ± 0.54	63.12 ± 1.44 60.08 ± 2.55	67.64 ± 0.74	67.28 ± 0.97	64.80 ± 0.85
P-WL P-WL-C P-WL-UC	79.34 ± 0.46 78.66 ± 0.32 78.50 ± 0.41	86.10 ± 1.37 90.51 ± 1.34 85.17 ± 0.29	85.34 ± 0.14 85.46 ± 0.16 85.62 ± 0.27	84.78 ± 0.15 84.96 ± 0.34 85.11 ± 0.30	75.31 ± 0.73 75.27 ± 0.38 75.86 ± 0.78	63.07 ± 1.68 64.02 ± 0.82 63.46 ± 1.58	67.30 ± 1.50 67.15 ± 1.09 67.02 ± 1.29	68.40 ± 1.17 68.57 ± 1.76 68.01 ± 1.04	64.47 ± 1.84 65.78 ± 1.22 65.44 ± 1.18

### **Graph representations**

Fundamental properties

- $^{\text{tr}}$  Two graphs  $\mathcal G$  and  $\mathcal G'$  can have a *different* number of vertices.
- $\hat{v}$  Hence, we require a vectorised representation  $f: \mathcal{G} \to \mathbb{R}^d$  of graphs.
- $\Rightarrow$  Such a representation f needs to be *permutation-invariant*.

### Now and then

#### Shallow approaches

- ☆ node2vec (encoder–decoder)
- ☆ Graph kernels (RKHS feature maps)
- ☆ Laplacian-based embeddings

### Deep approaches

- ☆ Graph convolutional networks
- ☆ Graph isomorphism networks
- Graph attention networks

The predominant paradigm in graph machine learning



- ☆ Operations remain local.
- ☆ Only require some aggregation function.
- Representations can be combined.

The predominant paradigm in graph machine learning



- ☆ Operations remain local.
- ☆ Only require some aggregation function.
- ☆ Representations can be combined.

The predominant paradigm in graph machine learning



- ☆ Operations remain local.
- © Only require some aggregation function.
- ☆ Representations can be combined.

The predominant paradigm in graph machine learning



- ☆ Operations remain local.
- ☆ Only require some aggregation function.
- ☆ Representations can be combined.

The predominant paradigm in graph machine learning



- ☆ Operations remain local.
- ☆ Only require some aggregation function.
- ☆ Representations can be combined.

The predominant paradigm in graph machine learning



- ☆ Operations remain local.
- ☆ Only require some aggregation function.
- Representations can be combined.

### Graph neural networks in a nutshell

$$\begin{split} a_v^{(k)} &:= \texttt{aggregate}^{(k)} \Big( \Big\{ h_u^{(k-1)} \mid u \in \mathcal{N}_{\mathcal{G}}(v) \Big\} \Big) \\ h_v^{(k)} &:= \texttt{combine}^{(k)} \Big( h_v^{(k-1)}, a_v^{(k)} \Big) \\ h_{\mathcal{G}} &:= \texttt{readout} \Big( \Big\{ h_v^{(K)} \mid v \in \texttt{vert}_{\mathcal{G}} \Big\} \Big) \end{split}$$

This terminology follows K. Xu, W. Hu, J. Leskovec and S. Jegelka, 'How Powerful are Graph Neural Networks?', *ICLR*, 2019.

### Example

Graph convolutional networks

For this architecture, combine is directly integrated into aggregate. In matrix form, we have

$$\mathbf{H}^{(\mathbf{k})} = \sigma \left( \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}} \mathbf{H}^{(\mathbf{k}-1)} \Theta^{(k-1)} \right)$$

with  $\mathbf{A} := \mathbf{A} + \mathbf{I}$  being the augmented adjacency matrix,  $\mathbf{D}$  its degree matrix, and  $\Theta^{(k-1)}$  a learnable weight matrix of dimensions  $d^{k-1} \times d^k$ .

### A topological layer for graph classification

Published as a conference paper at ICLR 2022

TOPOLOGICAL GRAPH NEURAL NETWORKS

Max Raw<sup>1,1,1</sup> Edward In Rosewer<sup>1,1</sup> Michael Max<sup>1,2</sup> You Messar<sup>2</sup> Radia Rivk<sup>1,1,1,1,1</sup> Kardes Regrand<sup>1,1,1</sup>

Department of Encyclement Advances and Englwarring TEE Zweich, 4000 Rand, Neckonstant Ref. Facility International Control of Constrained RAGE CARLER, KU Laveres, NWI Laveres, Refgion Industor of 24 for Reach, Histohen Studies, RCM Distoheng Genesary Tachakad Channesity of Nacida, NUCL Statistic, Common Tachakad Channesity of Nacida, NUCL Statistics, Common Tachakad Statistics, Nuclear Statistics, Nu

ABSTRACT

Ough mean strends (SDN), or a spaceful addition for is abling gaph beam of galaxy above means the observations in minore ubservations such a cybic-N garwaw TOCL , a need sport that incorporate galarit impedpate all and prove the term of the sport of the strength strength strength of the space of the sport of the strength strength strength strength strength strength incompletion with the same represents in items the Mathilan Lehnan apple incompletion with the same represents in items the Mathilan Lehnan apple incompletion with the same represents in the strength s

1 INTRODUCTION

Dopin as a reached proposation of messarily in the map domain, including bioinformatics are shown as the map down in the shown in the shown in the shown is the shown in the shown is the

To allow the two to properties  $q_{\rm eff}^{-1}$  again building figure (prover (DOG), that can be ready integrated as any GDM matching the properties of the second structure of the second structure

Our read-finations: We propose TOGL, as small layer based on TDA concepts that can be inteprinted into any CPNC. One hopes in differences and applied of lowering concenting topological distances of the state of the addity is not well with melli scate bargeness of the state of the state of the TOGC approve performance of averand GNN architectures when impeding in informations in strength TDM approximation of averand GNN architectures when impeding in informations is strength TDM approximate to the strengt GNN architectures when impeding in informations is a strength TDM approximate to the strengt GNN architectures when impeding in informations is a strength TDM and the strength of the strengt GNN architectures when impeding in informations is a strength TDM approximate to the strengt GNN architectures when impeding in informations is a strength TDM and the strength of the strengt GNN architectures when impeding in informations is a strength TDM architecture of the strengt GNN architecture of the strengt GNN architecture of the strength of the streng



Max Horn ♥@ExpectationMax



Yves Moreau



Edward De Brouwer



Karsten Borgwardt ♥@kmborgwardt



Michael Moor

M. Horn\*, E. De Brouwer\*, M. Moor, Y. Moreau, **B. Rieck**<sup>†</sup> and K. Borgwardt<sup>†</sup>, 'Topological Graph Neural Networks', *ICLR*, 2022

### **Motivation**

#### Status quo

- ☆ Graphs are topological objects.
- ☆ But GNNs are *incapable* of recognising certain topological structures!

### Challenge

What can we gain when imbuing them with knowledge about the topology?

### **Taking stock**

- ☆ Filtrations provide multi-scale topological features.
- ☆ Persistence diagrams serve as topological descriptors.

### Questions

- ☆ How to obtain 'good' filtrations?
- ☆ How to use persistence diagrams (i.e. multi-sets) in a differentiable setting?

## Topological graph neural networks

Overview



 $\hat{v}$  Use a node map  $\Phi \colon \mathbb{R}^d \to \mathbb{R}^k$  to create k different filtrations of the graph.

### Choosing $\Phi \, {\rm and} \, \Psi$

- $\Rightarrow$  The node map  $\Phi$  can be realised using a *neural network*.
- The coordinatisation function  $\Psi$  can be realised using *any* vectorisation of persistence diagrams (landscapes, images, ...), but we found a *differentiable coordinatisation function* to be most effective.<sup>2</sup>

<sup>2</sup>C. D. Hofer, F. Graf, **B. Rieck**, M. Niethammer and R. Kwitt, 'Graph Filtration Learning', *ICML*, 2020.

### **Expressivity of TOGL**

#### Theorem

TOGL (and persistent homology) is **more expressive** than WL[1], i.e. (i) if the WL[1] label sequences for two graphs G and G' diverge, there exists an injective filtration f such that the corresponding persistence diagrams  $D_0$  and  $D'_0$  are not equal, and (ii) there are graphs that WL[1] cannot distinguish but TOGL can!

### **Example graphs**





### Experiments

- ☆ Take existing GNN architecture.
- ☆ Replace one layer by TOGL.
- ☆ Measure predictive performance.

This strategy ensures that the number of parameters is approximately the same, thus facilitating a fair comparison!

### Synthetic data sets

Binary classification problem; generate same number of graphs for each of the classes. Use simple topological structures that are nevertheless challenging to detect with standard GNNs.



### Cycles data set

Weisfeiler-Lehman subtree features



These graphs cannot be distinguished based on their WL[1] information.

## Expressivity

Cycles data set



### Necklaces data set

Weisfeiler-Lehman subtree features



These graphs cannot be distinguished based on their WL[1] information (some of the graphs in the data set *can* be distinguished, though).

## Expressivity

Necklaces data set



### Classifying graphs/nodes based on structural features alone

Existing data sets tend to 'leak' information into node attributes, thus decreasing the utility of topological features. Hence, we replaced all node features by random ones.

	Node classification				
Метнор	DD	ENZYMES	MNIST	PROTEINS	Pattern
GCN-4 GCN-3-TOGL-1	$68.0 \pm 3.6$ <b>75.1 <math>\pm</math> 2.1</b>	$\begin{array}{c} 22.0 \pm 3.3 \\ \textbf{30.3} \pm \textbf{6.5} \end{array}$	$\begin{array}{c} 76.2 \pm 0.5 \\ \textbf{84.8} \pm \textbf{0.4} \end{array}$	$68.8 \pm 2.8$ <b>73.8 <math>\pm</math> 4.3</b>	$85.5 \pm 0.4$ 86.6 ± 0.1
GIN-4 GIN-3-TOGL-1	$75.6 \pm 2.8$ <b>76.2 <math>\pm</math> 2.4</b>	$21.3 \pm 6.5 \\ \textbf{23.7} \pm \textbf{6.9}$	$83.4 \pm 0.9$ 84.4 $\pm$ 1.1	$\begin{array}{c} {\bf 74.6 \pm 3.1} \\ {\bf 73.9 \pm 4.9} \end{array}$	$84.8 \pm 0.0$ <b>86.7 ± 0.1</b>
GAT-4 GAT-3-TOGL-1	$63.3 \pm 3.7$ <b>75.7 <math>\pm</math> 2.1</b>	$21.7 \pm 2.9 \\ \mathbf{23.5 \pm 6.1}$	$63.2 \pm 10.4$ <b>77.2 <math>\pm</math> 10.5</b>	$67.5 \pm 2.6$ <b>72.4 <math>\pm</math> 4.6</b>	$   \begin{array}{r} {\bf 73.1 \pm 1.9} \\ {59.6 \pm 3.3} \end{array} $

### Classifying benchmark data sets

While we improve baseline classification performance, the best performance is *not* driven by the availability of topological structures!

Graph classification								
Метнор	CIFAR-10	DD	ENZYMES	MNIST	PROTEINS-full	IMDB-B	REDDIT-B	CLUSTER
GATED-GCN-4 WL WL-OA	67.3 ± 0.3 —	72.9 ± 2.1 77.7 ± 2.0 <b>77.8 ± 1.2</b>	65.7 ± 4.9 54.3 ± 0.9 58.9 ± 0.9	97.3 ± 0.1 	<b>76.4 ± 2.9</b> 73.1 ± 0.5 73.5 ± 0.9	 71.2 ± 0.5 74.0 ± 0.7	 78.0 <u>+</u> 0.6 87.6 <u>+</u> 0.3	60.4 <u>+</u> 0.4 
GCN-4 GCN-3-TOGL-1	54.2 ± 1.5 61.7 ± 1.0 7.5	72.8 ± 4.1 73.2 ± 4.7 0.4	<b>65.8 ± 4.6</b> 53.0 ± 9.2 –12.8	90.0 ± 0.3 95.5 ± 0.2 5.5	76.1 ± 2.4 76.0 ± 3.9 -0.1	68.6 ± 4.9 72.0 ± 2.3 3.4	<b>92.8 ± 1.7</b> 89.4 ± 2.2 -3.4	57.0 ± 0.9 60.4 ± 0.2 3.4
GIN-4 GIN-3-TOGL-1	54.8 ± 1.4 61.3 ± 0.4 6.5	70.8 ± 3.8 75.2 ± 4.2 4.4	50.0 ± 12.3 43.8 ± 7.9 -6.2	96.1 ± 0.3 96.1 ± 0.1 0.0	72.3 ± 3.3 73.6 ± 4.8 1.3	72.8 ± 2.5 <b>74.2 ± 4.2</b> 1.4	81.7 ± 6.9 89.7 ± 2.5 8.0	58.5 ± 0.1 60.4 ± 0.2 1.9
GAT-4 GAT-3-TOGL-1	57.4 ± 0.6 63.9 ± 1.2 6.5	71.1 ± 3.1 73.7 ± 2.9 2.6	26.8 ± 4.1 51.5 ± 7.3 24.7	94.1 ± 0.3 95.9 ± 0.3 1.8	71.3 ± 5.4 75.2 ± 3.9 3.9	73.2 ± 4.1 70.8 ± 8.0 -2.4	44.2 ± 6.6 89.5 ± 8.7 45.3	56.6 ± 0.4 58.4 ± 3.7 1.8

### Conclusion

- 'If all you have is nails, everything looks like a hammer.'<sup>3</sup> Our data sets may actually stymie progress in GNN research because their classification does not necessarily require structural information.
- Nevertheless, higher-order structures (such as cliques) can be crucial in discerning between different graphs or data sets.
- ☆ Can we also learn sparse filtrations?

### Acknowledgements

My co-authors Edward, Karsten, Max, Michael, and Yves.

#### Software

https://github.com/aidos-lab/pytorch-topological
Looking for additional contributors!

<sup>3</sup>Credit: Mikael Vejdemo-Johannson