

Neural Sheaf Diffusion

A Topological Perspective on Heterophily and Oversmoothing in GNNs

Cristian Bodnar

Department of Computer Science
University of Cambridge

Based on joint work with Francesco Di Giovanni, Benjamin Chamberlain, Pietro Liò, Michael M. Bronstein

Recent Advances in Graph Machine Learning Workshop
Sorbonne Université
March 8, 2022

Motivation

Graph Convolutional Networks

Let $G = (V, E)$ be a graph with node feature matrix \mathbf{X} , adjacency matrix \mathbf{A} , degree matrix \mathbf{D} and normalised Laplacian Δ_0 . Consider the GCN [Kipf & Welling, 2017] equation:

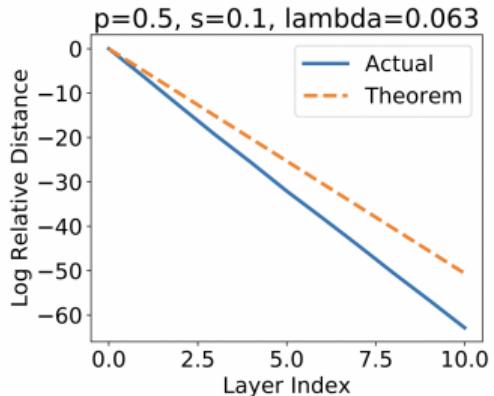
$$\begin{aligned} \text{GCN}(\mathbf{X}, \mathbf{A}) &:= \sigma\left((\mathbf{D} + \mathbf{I})^{-1/2}(\mathbf{A} + \mathbf{I})(\mathbf{D} + \mathbf{I})^{-1/2}\mathbf{X}\mathbf{W}\right) = \sigma\left(\tilde{\mathbf{D}}^{-1/2}\tilde{\mathbf{A}}\tilde{\mathbf{D}}^{-1/2}\mathbf{X}\mathbf{W}\right) \\ &= \sigma\left((\mathbf{I} - \tilde{\Delta}_0)\mathbf{X}\mathbf{W}\right) \end{aligned} \tag{1}$$

From Equation 1, it is clear the GCN is nothing but a non-linear, parametric and discretised diffusion equation:

$$\dot{\mathbf{X}}(t) = -\tilde{\Delta}_0\mathbf{X}(t) \rightsquigarrow \mathbf{X}(t+1) = \mathbf{X}(t) - \tilde{\Delta}_0\mathbf{X}(t) = (\mathbf{I} - \tilde{\Delta}_0)\mathbf{X}(t) \tag{2}$$

The Oversmoothing Problem

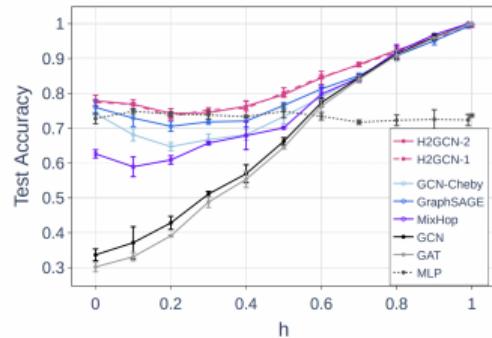
In Deep Graph Neural Networks (GNNs), it has been observed [Oono & Suzuki, 2020, Cai & Wang, 2020] that the features become progressively smoother with increased depth, resulting in a drop in the node classification performance. This is known as the “**oversmoothing problem**” of GNNs.



With more layers, GCN approaches a “smooth” subspace where all the node features are constant [Oono & Suzuki, 2020].

The Heterophily Problem

Numerous studies (e.g. [J. Zhu et al., 2020]) remarked that GNNs struggle in heterophilic settings (i.e. graphs where a node tends to be connected to nodes belonging to other classes). We call this the “**heterophily problem**” of GNNs.



The performance of GNNs is strongly correlated to the homophily level of a graph [J. Zhu et al., 2020].

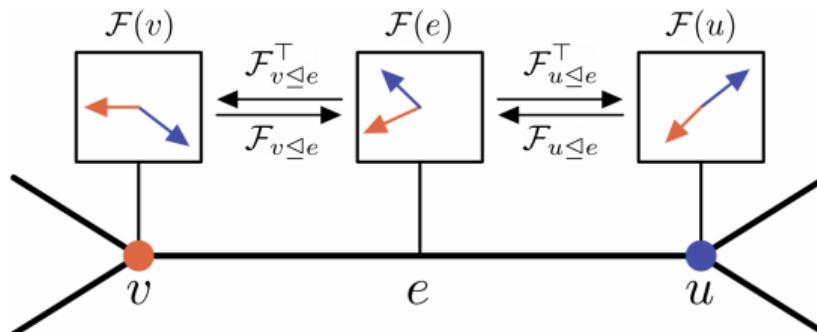
Cellular Sheaves

Sheaves

Definition

A *cellular sheaf* (G, \mathcal{F}) on an undirected graph $G = (V, E)$ consists of:

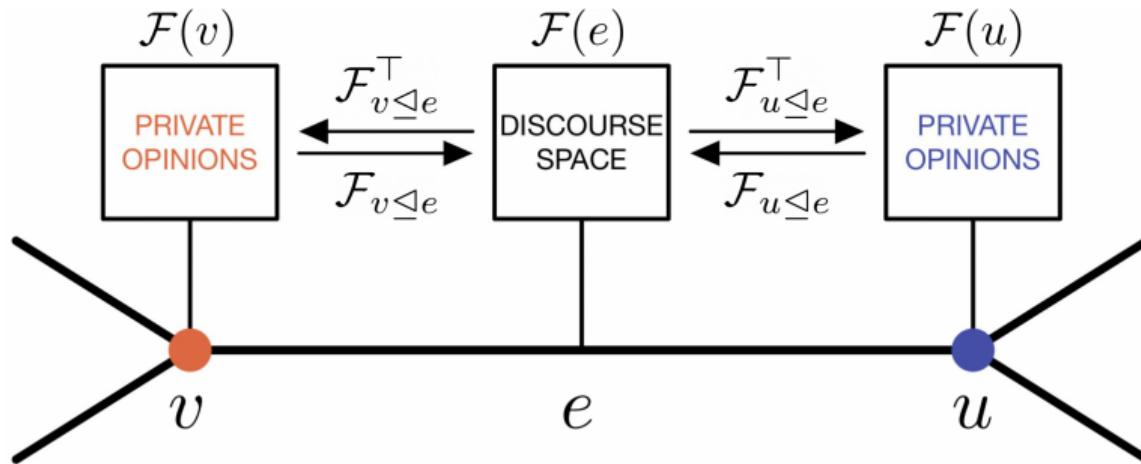
- A vector space $\mathcal{F}(v)$ for each $v \in V$.
- A vector space $\mathcal{F}(e)$ for each $e \in E$.
- A linear map $\mathcal{F}_{v \trianglelefteq e} : \mathcal{F}(v) \rightarrow \mathcal{F}(e)$ for each incident $v \trianglelefteq e$ node-edge pair.



The vector spaces are called *stalks* and the linear maps are also known as *restriction maps*.

Opinion Dynamics

From an opinion dynamics perspective [Hansen & Ghrist, 2020], if the nodes represent people, $\mathcal{F}(v)$ represents the space of private opinions of an individual and $\mathcal{F}_{v \leq e}$ encodes how this private opinion manifest in a discourse space $\mathcal{F}(e)$.



Cochains and the coboundary maps

We can form a space formed of all the spaces $\mathcal{F}(v)$ for all nodes v . Similarly, we can define a space formed of all the individual $\mathcal{F}(e)$.

Definition

For a sheaf (\mathcal{F}, G) we define the space of 0-cochains $C^0(G; \mathcal{F}) := \bigoplus_{v \in V} \mathcal{F}(v)$ and 1-cochains $C^1(G; \mathcal{F}) := \bigoplus_{e \in E} \mathcal{F}(e)$.

It is natural to define a linear *co-boundary map* δ between $C^0(G, \mathcal{F})$ and $C^1(G, \mathcal{F})$, which measures the disagreement between all nodes in the discourse space.

Definition

For some arbitrary choice of orientation for each edge $e = u \rightarrow v \in E$,
 $\delta: C^0(G, \mathcal{F}) \rightarrow C^1(G, \mathcal{F})$, $\delta(\mathbf{x})_e := \mathcal{F}_{v \triangleleft e} \mathbf{x}_v - \mathcal{F}_{u \triangleleft e} \mathbf{x}_u$.

The Sheaf Laplacian

We can use the co-boundary map to define a *sheaf Laplacian* [Hansen & Ghrist, 2019].

Definition

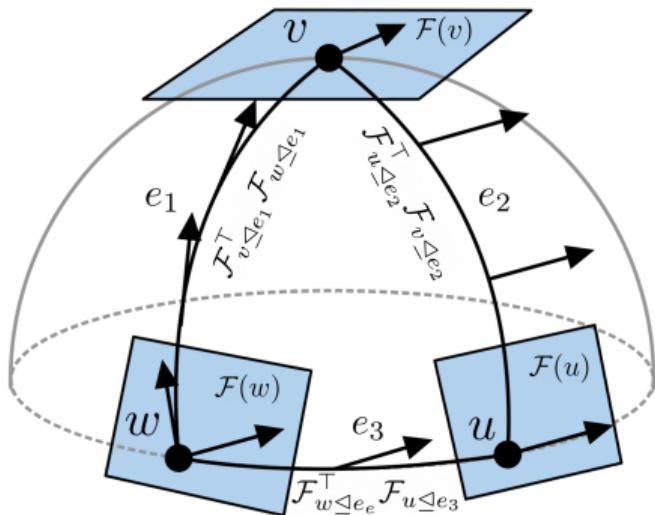
The sheaf Laplacian of a sheaf (G, \mathcal{F}) is a map $L_{\mathcal{F}} : C^0(G, \mathcal{F}) \rightarrow C^0(G, \mathcal{F})$ given by $L_{\mathcal{F}} := \delta^\top \delta$ or, equivalently, $L_{\mathcal{F}}(\mathbf{x})_v := \sum_{v, u \leq e} \mathcal{F}_{v \leq e}^\top (\mathcal{F}_{v \leq e} \mathbf{x}_v - \mathcal{F}_{u \leq e} \mathbf{x}_u)$

$$L_{\mathcal{F}} := \begin{bmatrix} v & \cdots & \cdots & \cdots & u \\ \vdots & & & & \vdots \\ \mathcal{F}_{v \leq e}^\top \mathcal{F}_{v \leq e} & \cdots & \cdots & \cdots & -\mathcal{F}_{v \leq e}^\top \mathcal{F}_{u \leq e} \\ \vdots & & \ddots & \ddots & \vdots \\ -\mathcal{F}_{u \leq e}^\top \mathcal{F}_{v \leq e} & \cdots & \cdots & \sum_{u \leq e} \mathcal{F}_{u \leq e}^\top \mathcal{F}_{u \leq e} \end{bmatrix}$$

The normalised Laplacian $\Delta_{\mathcal{F}} := D^{-1/2} L_{\mathcal{F}} D^{-1/2}$, where D is the block-diagonal of $L_{\mathcal{F}}$. Assuming all stalks have dimension d and the graph has n nodes, this is an $nd \times nd$ matrix, generalising the $n \times n$ (normalised) graph Laplacian matrix.

Discrete Vector Bundles

The sheaves (G, \mathcal{F}) with orthogonal restriction maps (i.e. $\mathcal{F}_{v \leq e} \in O(d)$) are called *discrete $O(d)$ -bundles* because they are a discrete equivalent of vector bundles.



Analogy between parallel transport on a sphere and transport on a discrete vector bundle. A tangent vector is moved from $\mathcal{F}(w) \rightarrow \mathcal{F}(v) \rightarrow \mathcal{F}(u)$ and back.

The Power of Sheaf Diffusion

Sheaf Diffusion

We can represent the d -dimensional vector-features of the nodes as a matrix $\mathbf{X} \in \mathbb{R}^{(nd) \times f}$ with f feature channels, whose columns are vectors in $C^0(G; \mathcal{F})$.

We are interested in the asymptotic behaviour of spatially discretised *sheaf diffusion* process governed by the PDE:

$$\mathbf{X}(0) = \mathbf{X}, \quad \dot{\mathbf{X}}(t) = -\Delta_{\mathcal{F}} \mathbf{X}(t)$$

Theorem (Hodge Theorem [Hansen & Ghrist, 2020])

In the infinite-time limit, the features converge to the projection of $\mathbf{X}(0)$ into $\ker(\Delta_{\mathcal{F}})$, which is the subspace of signals \mathbf{x} such that $\mathcal{F}_{v \trianglelefteq e} D_v^{-1/2} \mathbf{x}_v = \mathcal{F}_{u \trianglelefteq e} D_u^{-1/2} \mathbf{x}_u$ for all $e \in E$.

Intuition. All the private opinions converge towards a configuration where everyone agrees with all their neighbours in the discourse space (up to a $D^{-1/2}$ normalisation).

The Separation Power of Sheaf Diffusion

Definition

A hypothesis class of sheaves with d -dimensional stalks \mathcal{H}^d has linear separation power over a set of graphs \mathcal{G} if for any labelled graph $G = (V, E) \in \mathcal{G}$, there is a sheaf $(\mathcal{F}, G) \in \mathcal{H}^d$ that can linearly separate the classes of G in the diffusion time-limit for a dense subset $\mathcal{X}_{\mathcal{F}} \subset \mathbb{R}^{nd \times f}$ of initial conditions.

How is this framework related to oversmoothing and heterophily?

1. Linear separation power over $\mathcal{G} \iff$ No oversmoothing over \mathcal{G} .
2. Does \mathcal{G} contain (many) heterophilic graphs? \iff Robust to heterophilic settings.

The Curse of Symmetric Transport

Definition

Consider the class of sheaves with symmetric and invertible transport maps and d -dimensional stalks: $\mathcal{H}_{\text{sym}}^d := \{(\mathcal{F}, G) \mid \mathcal{F}_{v \leq e} = \mathcal{F}_{u \leq e}, \det(\mathcal{F}_{v \leq e}) \neq 0\}$

Proposition

Let \mathcal{G} be the set of connected graphs $G = (V, E)$ with two classes $A, B \subset V$ such that for each $v \in A$, there exists $u \in A$ and an edge $(v, u) \in E$. Then $\mathcal{H}_{\text{sym}}^1$ has linear separation power over \mathcal{G} .

Proposition

Let \mathcal{G} be the set of connected bipartite graphs $G = (A, B, E)$, with partitions A, B forming two classes and $|A| = |B|$. Then $\mathcal{H}_{\text{sym}}^1$ cannot linearly separate any graph in \mathcal{G} for any initial conditions $\mathbf{X}(0) \in \mathbb{R}^{n \times f}$.

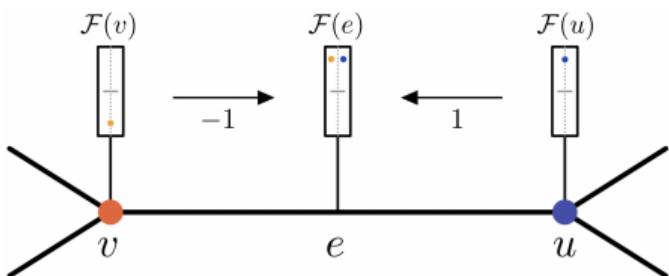
Lying is good

Definition

Dropping the equality, $\mathcal{H}_*^d := \{(\mathcal{F}, G) \mid \det(\mathcal{F}_{v \trianglelefteq e}) \neq 0\}$

Proposition

Let \mathcal{G} contain all the connected graphs with two classes. Then, \mathcal{H}_*^1 has linear separation power over \mathcal{G} .



Opposite signs lead to opinion polarisation.

This explains why a recent body of work [Yan et al., 2021, Chien et al., 2021] has found negatively-weighted edges to help in heterophilic settings. If $\mathcal{F}_{v \trianglelefteq e} \mathcal{F}_{u \trianglelefteq e} < 0$, then $\mathcal{F}_{v \trianglelefteq e} \neq \mathcal{F}_{u \trianglelefteq e}$.

The blessing of dimensionality

Even with all this additional flexibility, dimension $d = 1$ still has a major limitation.

Proposition

Let G be a connected graph with $C \geq 3$ classes. \mathcal{H}_*^1 cannot separate any $\mathbf{X}(0) \in \mathbb{R}^{n \times f}$.

With sufficient stalk dimension (i.e. width), this can be fixed.

Definition

Consider the class of sheaves with diagonal invertible maps and d -dimensional stalks:

$$\mathcal{H}_{\text{diag}}^d := \{(\mathcal{F}, G) \mid \mathcal{F}_{v \trianglelefteq e} = \text{invertible diagonal matrix}\}$$

Proposition

Let \mathcal{G} be the set of connected graphs with nodes belonging to $C \geq 3$ classes. Then for $d \geq C$, $\mathcal{H}_{\text{diag}}^d$ has linear separation power over \mathcal{G} .

Beyond diagonal maps

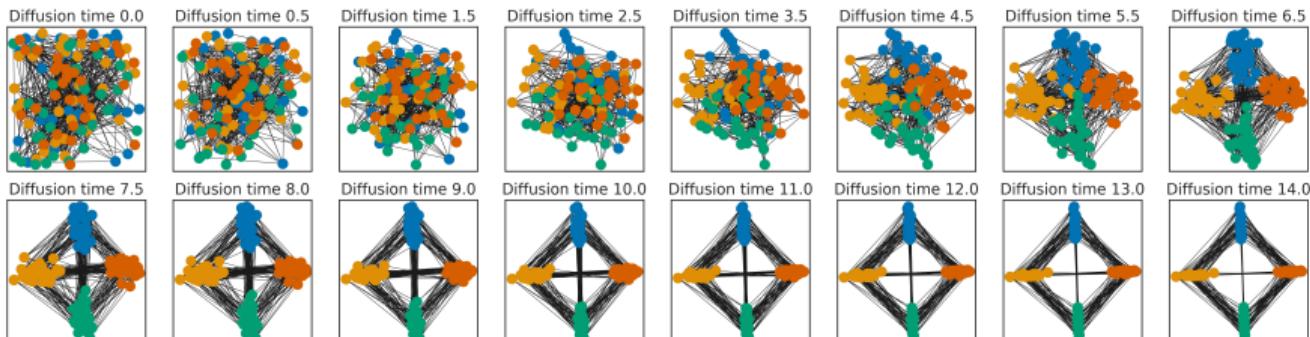
Diagonal restriction maps are extremely simple. Can we do better?

Definition

The class of discrete $O(d)$ -bundles $\mathcal{H}_{\text{orth}}^d := \{(\mathcal{F}, G) \mid \mathcal{F}_{v \leq e} \in O(d)\}$

Theorem

Let \mathcal{G} be the class of connected graphs with $C \leq 2d$ classes. Then, for all $d \in \{2, 4\}$, $\mathcal{H}_{\text{orth}}^d$ has linear separation power over \mathcal{G} .



The Power of Sheaf Diffusion: Review

In summary, we showed that:

1. Sheaf diffusion on a trivial sheaf with symmetric (scalar) restriction maps (as implicitly used in standard graph convolutions) has linear separation power only in certain (homophilic) settings.
2. Dealing with heterophilic data and oversmoothing requires non-symmetric restriction maps.
3. Higher-dimensional stalks and more complex restriction maps lead to stronger separation power.

Learning Sheaves

Learning Sheaves

In practice, the ground truth sheaf for a task is unknown, so we aim to learn it from data. We consider the following diffusion-type equation:

$$\dot{\mathbf{X}}(t) = -\sigma \left(\Delta_{\mathcal{F}(t)} (\mathbf{I}_n \otimes \mathbf{W}_1) \mathbf{X}(t) \mathbf{W}_2 \right), \quad (3)$$

Crucially, the sheaf evolves over time as a function of the data $(G, \mathcal{F}(t)) = g(G, \mathbf{X}(t); \theta)$.

We also consider a discrete version of this equation, with different weights at each layer t .

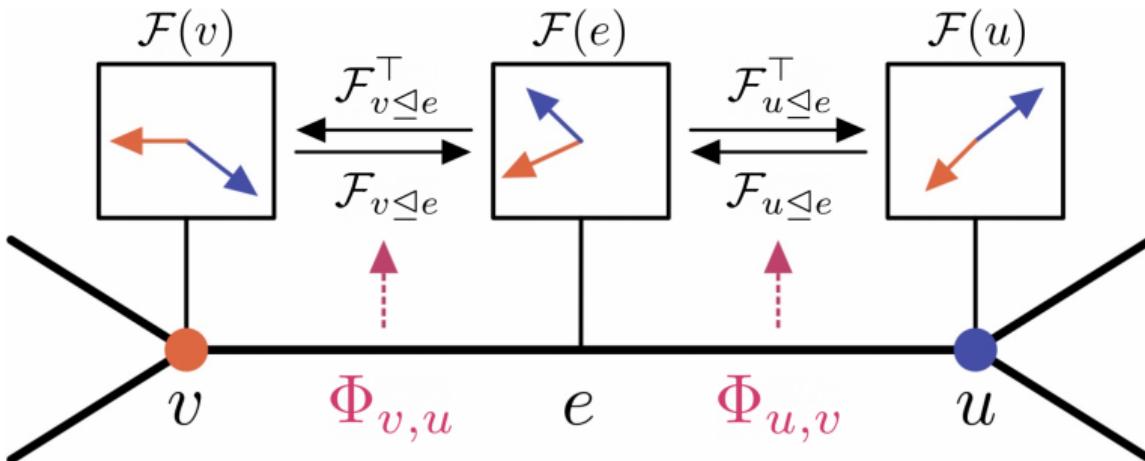
$$\mathbf{X}_{t+1} = \mathbf{X}_t - \sigma \left(\Delta_{\mathcal{F}(t)} (\mathbf{I} \otimes \mathbf{W}_1^t) \mathbf{X}_t \mathbf{W}_2^t \right), \quad (4)$$

Learning the restriction maps

Each $d \times d$ matrix $\mathcal{F}_{v \trianglelefteq e}$ is learned via a parametric function $\Phi : \mathbb{R}^{d \times 2} \rightarrow \mathbb{R}^{d \times d}$:

$$\mathcal{F}_{v \trianglelefteq e:=(v,u)} = \Phi(\mathbf{x}_v, \mathbf{x}_u) \quad (5)$$

The restriction map can be **diagonal**, **orthogonal**, or a **general** matrix.



The restriction maps are learned from data.

Computational Complexity

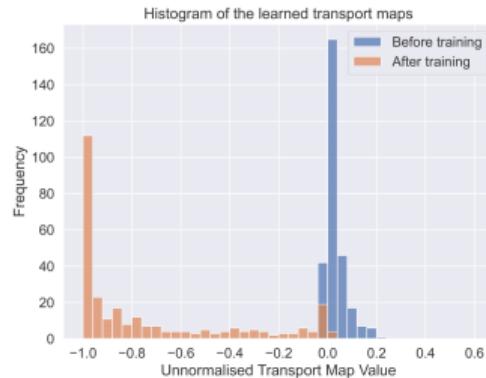
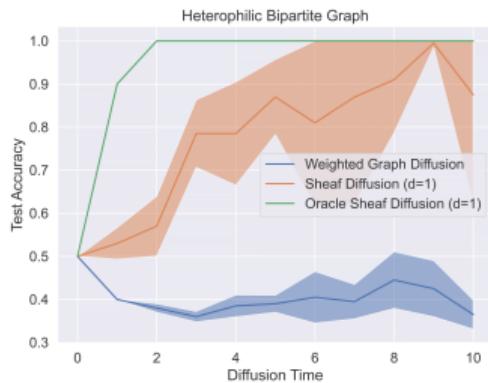
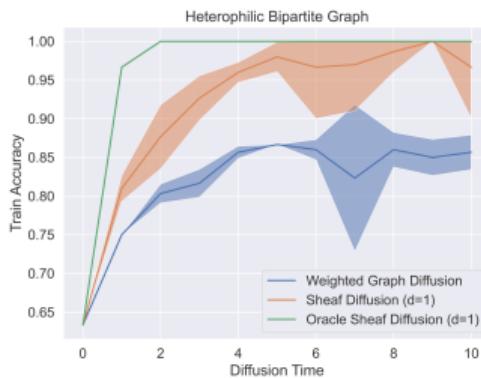
Let n, m be the number of nodes and edges of the graph and consider that all stalks have dimension d . The computational complexity varies with the type of sheaf we learn:

- **Diagonal:** Matrix multiplication becomes scalar multiplication, so complexity is $\mathcal{O}(d(n + m))$.
- **Orthogonal / General:** Because matrix multiplication is needed to construct the transport maps: $\mathcal{O}(d^3(n + m))$. However, at inference time, one can use a pre-computed Laplacian and complexity becomes $\mathcal{O}(d^2(n + m))$.
- **Other:** One can flexibly adjust between these two extremes by making the transport maps more sparse (e.g. block-diagonal).

Results

Synthetic Experiment: Opinion Polarisation

We have a bipartite graph with equally sized partitions that we try to distinguish. $\mathbf{X}(0)$ is not linearly separable. We use a simple sheaf diffusion process with a learned sheaf Laplacian (i.e. no weights and non-linearities)



Training (*Left*) and Testing (*Middle*) accuracy as a function of diffusion time. Learned sheaf Laplacian for $t \gg 0$. (*Right*)

Real-World Evaluation

We evaluate on multiple node-classifications tasks with various degrees of homophily [Rozemberczki et al, 2019, Pei et al, 2019].

	Texas	Wisconsin	Film	Squirrel	Chameleon	Cornell	Citeseer	Pubmed	Cora
Hom level	0.11	0.21	0.22	0.22	0.23	0.30	0.74	0.80	0.81
#Nodes	183	251	7,600	5,201	2,277	183	3,327	18,717	2,708
#Edges	295	466	26,752	198,493	31,421	280	4,676	44,327	5,278
#Classes	5	5	5	5	5	5	7	3	6
Diag-SD	85.67 \pm 6.95	88.63 \pm 2.75	37.79 \pm 1.01	54.78 \pm 1.81	68.68 \pm 1.73	86.49 \pm 7.35	77.14 \pm 1.85	89.42 \pm 0.43	87.14 \pm 1.06
O(d)-SD	85.95 \pm 5.51	89.41 \pm 4.74	37.81 \pm 1.15	56.34 \pm 1.32	68.04 \pm 1.58	84.86 \pm 4.71	76.70 \pm 1.57	89.49 \pm 0.40	86.90 \pm 1.13
Gen-SD	82.97 \pm 5.13	89.21 \pm 3.84	37.80 \pm 1.22	53.17 \pm 1.31	67.93 \pm 1.58	85.68 \pm 6.51	76.32 \pm 1.65	89.33 \pm 0.35	87.30 \pm 1.15
GGCN	84.86 \pm 4.55	86.86 \pm 3.29	37.54 \pm 1.56	55.17 \pm 1.58	71.14 \pm 1.84	85.68 \pm 6.63	77.14 \pm 1.45	89.15 \pm 0.37	87.95 \pm 1.05
H2GCN	84.86 \pm 7.23	87.65 \pm 4.98	35.70 \pm 1.00	36.48 \pm 1.86	60.11 \pm 2.15	82.70 \pm 5.28	77.11 \pm 1.57	89.49 \pm 0.38	87.87 \pm 1.20
GPRGNN	78.38 \pm 4.36	82.94 \pm 4.21	34.63 \pm 1.22	31.61 \pm 1.24	46.58 \pm 1.71	80.27 \pm 8.11	77.13 \pm 1.67	87.54 \pm 0.38	87.95 \pm 1.18
FAGCN	82.43 \pm 6.89	82.94 \pm 7.95	34.87 \pm 1.25	42.59 \pm 0.79	55.22 \pm 3.19	79.19 \pm 9.79	N/A	N/A	N/A
MixHop	77.84 \pm 7.73	75.88 \pm 4.90	32.22 \pm 2.34	43.80 \pm 1.48	60.50 \pm 2.53	73.51 \pm 6.34	76.26 \pm 1.33	85.31 \pm 0.61	87.61 \pm 0.85
GCNII	77.57 \pm 3.83	80.39 \pm 3.40	37.44 \pm 1.30	38.47 \pm 1.58	63.86 \pm 3.04	77.86 \pm 3.79	77.33 \pm 1.48	90.15 \pm 0.43	88.37 \pm 1.25
Geom-GCN	66.76 \pm 2.72	64.51 \pm 3.66	31.59 \pm 1.15	38.15 \pm 0.92	60.00 \pm 2.81	60.54 \pm 3.67	78.02 \pm 1.15	89.95 \pm 0.47	85.35 \pm 1.57
PairNorm	60.27 \pm 4.34	48.43 \pm 6.14	27.40 \pm 1.24	50.44 \pm 2.04	62.74 \pm 2.82	58.92 \pm 3.15	73.59 \pm 1.47	87.53 \pm 0.44	85.79 \pm 1.01
GraphSAGE	82.43 \pm 6.14	81.18 \pm 5.56	34.23 \pm 0.99	41.61 \pm 0.74	58.73 \pm 1.68	75.95 \pm 5.01	76.04 \pm 1.30	88.45 \pm 0.50	86.90 \pm 1.04
GCN	55.14 \pm 5.16	51.76 \pm 3.06	27.32 \pm 1.10	53.43 \pm 2.01	64.82 \pm 2.24	60.54 \pm 5.30	76.50 \pm 1.36	88.42 \pm 0.50	86.98 \pm 1.27
GAT	52.16 \pm 6.63	49.41 \pm 4.09	27.44 \pm 0.89	40.72 \pm 1.55	60.26 \pm 2.50	61.89 \pm 5.05	76.55 \pm 1.23	87.30 \pm 1.10	86.33 \pm 0.48
MLP	80.81 \pm 4.75	85.29 \pm 3.31	36.53 \pm 0.70	28.77 \pm 1.56	46.21 \pm 2.99	81.89 \pm 6.40	74.02 \pm 1.90	75.69 \pm 2.00	87.16 \pm 0.37
Cont Diag-SD	82.97 \pm 4.37	86.47 \pm 2.55	36.85 \pm 1.21	38.17 \pm 9.29	62.06 \pm 3.84	80.00 \pm 6.07	76.56 \pm 1.19	89.47 \pm 0.42	86.88 \pm 1.21
Cont O(d)-SD	82.43 \pm 5.95	84.50 \pm 4.34	36.39 \pm 1.37	40.40 \pm 2.01	63.18 \pm 1.69	72.16 \pm 10.40	75.19 \pm 1.67	89.12 \pm 0.30	86.70 \pm 1.24
Cont Gen-SD	83.78 \pm 6.62	85.29 \pm 3.31	37.28 \pm 0.74	52.57 \pm 2.76	66.40 \pm 2.28	84.60 \pm 4.69	77.54 \pm 1.72	89.67 \pm 0.40	87.45 \pm 0.99
BLEND	83.24 \pm 4.65	84.12 \pm 3.56	35.63 \pm 0.89	43.06 \pm 1.39	60.11 \pm 2.09	85.95 \pm 6.82	76.63 \pm 1.60	89.24 \pm 0.42	88.09 \pm 1.22
GRAND	75.68 \pm 7.25	79.41 \pm 3.64	35.62 \pm 1.01	40.05 \pm 1.50	54.67 \pm 2.54	82.16 \pm 7.09	76.46 \pm 1.77	89.02 \pm 0.51	87.36 \pm 0.96
CGNN	71.35 \pm 4.05	74.31 \pm 7.26	35.95 \pm 0.86	29.24 \pm 1.09	46.89 \pm 1.66	66.22 \pm 7.69	76.91 \pm 1.81	87.70 \pm 0.49	87.10 \pm 1.35

References

-  Cristian Bodnar, Francesco Di Giovanni, Benjamin Paul Chamberlain, Pietro Liò, Michael M. Bronstein (2022)
Neural Sheaf Diffusion: A Topological Perspective on Heterophily and Oversmoothing in GNNs
Preprint 2022
-  Thomas N. Kipf, Max Welling
Semi-Supervised Classification with Graph Convolutional Networks
ICLR 2017
-  Kenta Oono, Taiji Suzuki (2020)
Graph Neural Networks Exponentially Lose Expressive Power for Node Classification
ICLR 2020
-  Jiong Zhu, Yujun Yan, Lingxiao Zhao, Mark Heimann, Leman Akoglu, Danai Koutra
Beyond Homophily in Graph Neural Networks: Current Limitations and Effective Designs
NeurIPS 2020

References



Jakob Hansen, Robert Ghrist

Opinion Dynamics on Discourse Sheaves

SIAM J. Appl. Math., 81(5), 2033–2060



Jakob Hansen, Robert Ghrist

Toward a Spectral Theory of Cellular Sheaves

Journal of Applied and Computational Topology volume 3, pages 315–358 (2019)



Yujun Yan, Milad Hashemi, Kevin Swersky, Yaoqing Yang, Danai Koutra

Two Sides of the Same Coin: Heterophily and Oversmoothing in Graph Convolutional Neural Networks

Preprint (2021)



Eli Chien, Jianhao Peng, Pan Li, Olgica Milenkovic

Adaptive Universal Generalized PageRank Graph Neural Network

ICLR 2021

References



Chen Cai, Yusu Wang

A Note on Over-Smoothing for Graph Neural Networks

ICML 2020 Graph Representation Learning Workshop



Jakob Hansen, Thomas Gebhart

Sheaf Neural Networks

NeurIPS 2020 Workshop on TDA and Beyond



Benedek Rozemberczki, Carl Allen, Rik Sarkar

Multi-scale Attributed Node Embedding

Journal of Complex Networks



Hongbin Pei, Bingzhe Wei, Kevin Chen-Chuan Chang, Yu Lei, Bo Yang

Geom-GCN: Geometric Graph Convolutional Networks

ICLR 2020

The End