

# The Effects of Activation Functions on the Over-smoothing Issue of Graph Convolutional Networks

Bao Wang

Department of Mathematics  
Scientific Computing and Imaging Institute  
The University of Utah

CBMS Conference: Deep Learning and Numerical PDEs  
Morgan State University



---

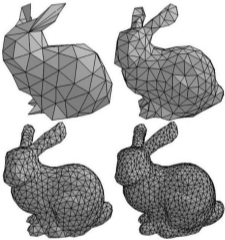
Shih-Hsin Wang\*, Justin Baker\*, Cory Hauck, and Bao Wang, The Effects of Activation Functions on the Over-smoothing Issue of Graph Convolutional Networks, submitted.

# Learning Non-Euclidean Data?

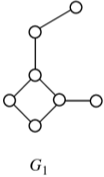
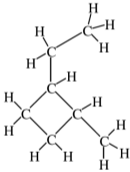
- Graph is a flexible structure to represent non-Euclidean data.



**Social Graph**  
(Facebook, Wikipedia)



**3D Mesh**



**Molecular Graph**

## Graph convolutional networks

- Let  $G = (V, E)$  be an undirected graph where  $V = \{v_i\}_{i=1}^n$  is the set of nodes and  $E$  is the set of edges.
- Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  be the adjacency matrix of  $G$ .
- Let  $\mathbf{G} := (\mathbf{D} + \mathbf{I})^{-\frac{1}{2}}(\mathbf{I} + \mathbf{A})(\mathbf{D} + \mathbf{I})^{-\frac{1}{2}} = \tilde{\mathbf{D}}^{-\frac{1}{2}}\tilde{\mathbf{A}}\tilde{\mathbf{D}}^{-\frac{1}{2}}$  be the (augmented) normalized adjacency matrix.
- Graph convolutional layer (GCL):

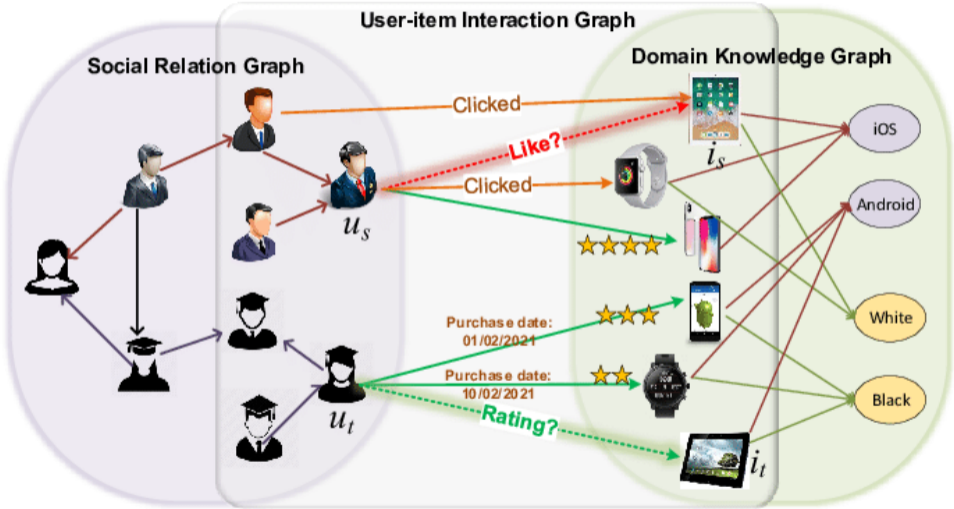
$$\mathbf{H}^l = \sigma(\mathbf{W}^l \mathbf{H}^{l-1} \mathbf{G}),$$

where  $\sigma$  is the activation function,  $\mathbf{W}^l \in \mathbb{R}^{d \times d}$  is a learnable weight matrix, and  $\mathbf{H}^0 := [\mathbf{h}_1, \dots, \mathbf{h}_n] \in \mathbb{R}^{d \times n}$  with  $\mathbf{h}_i$  being the  $i^{\text{th}}$  node feature. [A message-passing scheme rather than exact convolution.](#)

## Graph learning tasks

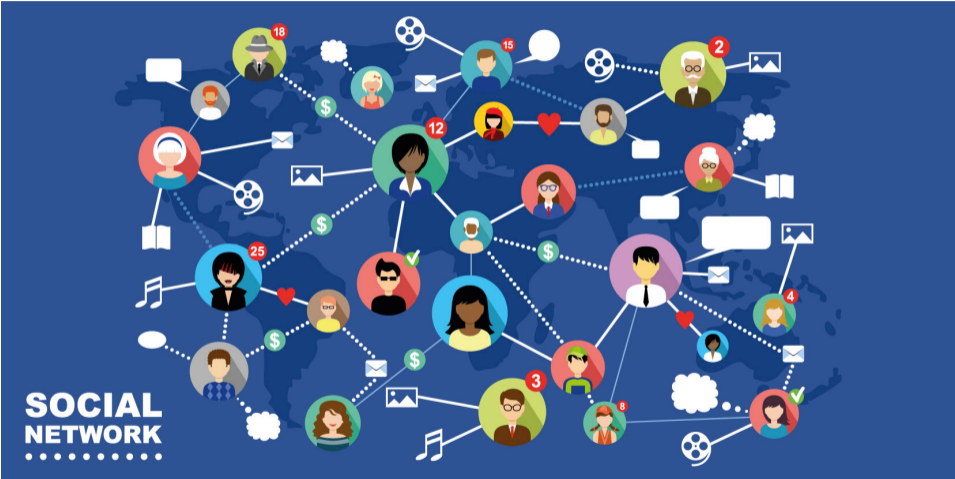
- Node classification
- Link prediction
- Graph classification and generation

Applications: Recommender system



Link prediction

Applications: Social network



Node classification

## Over-smoothing of GNN

- All eigenvalues of  $\mathbf{G}$  lie in the interval  $(-1, 1]$ .
- $\mathbf{H}^l = \mathbf{W}^l \mathbf{H}^{l-1} \mathbf{G}$ , i.e.,  $\text{vec}(\mathbf{H}^l) = \mathbf{G}^\top \otimes \mathbf{W}^l \text{vec}(\mathbf{H}^{l-1})$ , can be considered as a low-pass filter, indicating that each GCL “smooths” node features.
- As the GCN architecture gets deep, all nodes’ representation – within each connected component – will become “indistinguishable”, which is referred to as *over-smoothing*.
- Learning long-range dependencies (“long-range interaction”) is hard.



## Existing Theory

## Mathematical characterization of the over-smoothing - I (Oono & Suzuki, ICLR, 2019.)

- Distance of the node features  $\mathbf{H}^l$  to the eigenspace  $\mathcal{M}$  – the eigenspace corresponding to the largest eigenvalue of  $\mathbf{G}$  – goes to zero.

> Suppose the graph  $G$  has  $m$  connected components, i.e. we can decompose  $V = \bigcup_{i=1}^m V_i$ . Let  $\mathbf{u}_i = (\mathbf{1}_{\{k \in V_i\}})_{1 \leq k \leq n}$  be the indicator vector of the  $i^{\text{th}}$  component  $V_i$ .

> The nonnegative vectors  $\{\tilde{\mathbf{D}}^{\frac{1}{2}} \mathbf{u}_i / \|\tilde{\mathbf{D}}^{\frac{1}{2}} \mathbf{u}_i\|\}_{1 \leq i \leq m}$  form an orthonormal basis of  $\mathcal{M}$ .

- Let  $\mathbb{R}^d \otimes \mathcal{M}$  be the subspace of  $\mathbb{R}^{d \times n}$  consisting of the sum  $\sum_{i=1}^m \mathbf{w}_i \otimes \mathbf{e}_i$  where  $\mathbf{w}_i \in \mathbb{R}^d$  and  $\{\mathbf{e}_i\}_{i=1}^m$  is an orthonormal basis of the eigenspace  $\mathcal{M}$ . Then the distance of  $\mathbf{H}^l$  to  $\mathcal{M}$  is

$$\|\mathbf{H}^l\|_{\mathcal{M}^\perp} := \inf_{\mathbf{Y} \in \mathbb{R}^d \otimes \mathcal{M}} \|\mathbf{H}^l - \mathbf{Y}\|_F = \left\| \mathbf{H}^l - \sum_{i=1}^m \mathbf{H}^l \mathbf{e}_i \mathbf{e}_i^\top \right\|_F.$$

- $\|\mathbf{H}^l\|_{\mathcal{M}^\perp} \leq s_l \lambda \|\mathbf{H}^{l-1}\|_{\mathcal{M}^\perp}$  when  $\sigma$  is ReLU. Here,  $\lambda = \max\{|\lambda_i| \mid \lambda_i < 1\}$  is the second largest magnitude of  $\mathbf{G}$ 's eigenvalues, and  $s_l$  is the largest singular value of  $\mathbf{W}^l$ .

## Effects of ReLU

- $\|\sigma(\mathbf{Z})\|_{\mathcal{M}^\perp} \leq \|\mathbf{Z}\|_{\mathcal{M}^\perp}$  for any matrix  $\mathbf{Z}$  when  $\sigma$  is ReLU, i.e. ReLU reduces the distance to eigenspace  $\mathcal{M}$ . – [Oono & Suzuki, ICLR, 2019](#)

- Dirichlet energy of node features:

$$\|\mathbf{H}\|_E^2 := \text{Trace}(\mathbf{H}\tilde{\Delta}\mathbf{H}^\top),$$

where  $\tilde{\Delta} = \mathbf{I} - \mathbf{G}$  is the (augmented) normalized Laplacian.

- $\|\mathbf{H}'\|_E \leq s_l \lambda \|\mathbf{H}'^{-1}\|_E$  when  $\sigma$  is ReLU or leaky ReLU.

## Effects of activation function: Existing theory

- $\|\sigma(\mathbf{Z})\|_{\mathcal{M}^\perp} \leq \|\mathbf{Z}\|_{\mathcal{M}^\perp}$  for any matrix  $\mathbf{Z}$  when  $\sigma$  is ReLU, i.e. ReLU reduces the distance to eigenspace  $\mathcal{M}$ . – [Oono & Suzuki, ICLR, 2019](#)
  - $\|\sigma(\mathbf{Z})\|_E \leq \|\mathbf{Z}\|_E$  for any matrix  $\mathbf{Z}$  when  $\sigma$  is ReLU or leaky ReLU. – [Cai & Wang, arXiv:2006.13318, 2020](#)
  - $\|\mathbf{H}\|_{\mathcal{M}^\perp}$  and  $\|\mathbf{H}\|_E$  are two equivalent seminorms, i.e. there exist two constants  $\alpha, \beta > 0$  s.t.  $\alpha\|\mathbf{H}\|_{\mathcal{M}^\perp} \leq \|\mathbf{H}\|_E \leq \beta\|\mathbf{H}\|_{\mathcal{M}^\perp}$  for any  $\mathbf{H} \in \mathbb{R}^{d \times n}$ .
- >  $\|\sigma(\mathbf{Z})\|_{\mathcal{M}^\perp} \leq \|\mathbf{Z}\|_{\mathcal{M}^\perp}$ , when  $\sigma$  is ReLU or leaky ReLU.

## Bottlenecks of the existing theory

- Existing smoothness notions – distant to  $\mathcal{M}$  and Dirichlet energy of node features – do not take the magnitude of feature vectors into account and they are not scaling free. Multiplying feature vectors by a constant will result in corresponding changes in their distance to  $\mathcal{M}$  and their Dirichlet energy but do not affect graph node classification.
- Existing theory do not reveal a mechanism to control the smoothness of the learned node features when taking the activation functions into consideration.

# Geometry Underlying the Input & Output of ReLU and Leaky ReLU

## Geometric characterization of the effect of ReLU

- We have the decomposition  $\mathbf{H} = \mathbf{H}_{\mathcal{M}} + \mathbf{H}_{\mathcal{M}^\perp}$  for any matrix  $\mathbf{H} := [\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_n] \in \mathbb{R}^{d \times n}$

$$\mathbf{H}_{\mathcal{M}} = \sum_{i=1}^m \mathbf{H} \mathbf{e}_i \mathbf{e}_i^\top, \quad \text{and} \quad \mathbf{H}_{\mathcal{M}^\perp} = \sum_{i=m+1}^n \mathbf{H} \mathbf{e}_i \mathbf{e}_i^\top.$$

- Let  $\mathbf{Z} \in \mathbb{R}^{d \times n}$  be an arbitrary matrix and  $\mathbf{H} = \sigma(\mathbf{Z})$  with  $\sigma(x) = \max\{0, x\}$  being ReLU.
- **Proposition 1.** For any  $\mathbf{Z} = \mathbf{Z}_{\mathcal{M}} + \mathbf{Z}_{\mathcal{M}^\perp} \in \mathbb{R}^{d \times n}$ , let  $\mathbf{H} = \sigma(\mathbf{Z}) = \mathbf{H}_{\mathcal{M}} + \mathbf{H}_{\mathcal{M}^\perp}$  with  $\sigma$  being ReLU, then  $\mathbf{H}_{\mathcal{M}^\perp}$  lies on the high dimensional sphere centered at  $\mathbf{Z}_{\mathcal{M}^\perp}/2$  with the radius

$$r := \left( \|\mathbf{Z}_{\mathcal{M}^\perp}/2\|_F^2 - \langle \mathbf{Z}_{\mathcal{M}}^+, \mathbf{Z}_{\mathcal{M}}^- \rangle_F \right)^{1/2}.$$

In particular,  $\mathbf{H}_{\mathcal{M}^\perp}$  lies inside the ball centered at  $\mathbf{Z}_{\mathcal{M}^\perp}/2$  with radius  $\|\mathbf{Z}_{\mathcal{M}^\perp}/2\|_F$  and hence we have  $\|\mathbf{H}_{\mathcal{M}^\perp}\|_F \leq \|\mathbf{Z}_{\mathcal{M}^\perp}\|_F$ . **[Reduced distance to  $\mathcal{M}$ !]**

- $\mathbf{Z}^+ = \max(\mathbf{Z}, 0)$  and  $\mathbf{Z}^- = \max(-\mathbf{Z}, 0)$ .



## Geometric characterization of the effect of leaky ReLU

- Let  $\mathbf{Z} \in \mathbb{R}^{d \times n}$  be an arbitrary matrix and  $\mathbf{H} = \sigma_a(\mathbf{Z})$  with  $\sigma_a$  being leaky ReLU:

$$\sigma_a(x) = \begin{cases} x & \text{if } x \geq 0, \\ ax & \text{otherwise,} \end{cases}$$

where  $0 < a < 1$  is a positive scalar.

- Proposition 2.** For any  $\mathbf{Z} = \mathbf{Z}_{\mathcal{M}} + \mathbf{Z}_{\mathcal{M}^\perp}$ , let  $\mathbf{H} = \sigma_a(\mathbf{Z}) = \mathbf{H}_{\mathcal{M}} + \mathbf{H}_{\mathcal{M}^\perp}$  with  $\sigma_a$  being leaky ReLU, then  $\mathbf{H}_{\mathcal{M}^\perp}$  lies on the high dimensional sphere centered at  $(1+a)\mathbf{Z}_{\mathcal{M}^\perp}/2$  with radius

$$r_a := \left( \|(1-a)\mathbf{Z}_{\mathcal{M}^\perp}/2\|_F^2 - (1-a)^2 \langle \mathbf{Z}_{\mathcal{M}}^+, \mathbf{Z}_{\mathcal{M}}^- \rangle_F \right)^{1/2}.$$

In particular,  $\mathbf{H}_{\mathcal{M}^\perp}$  lies inside the high-dimensional ball centered at  $(1+a)\mathbf{Z}_{\mathcal{M}^\perp}/2$  with radius  $\|(1-a)\mathbf{Z}_{\mathcal{M}^\perp}/2\|_F$  and hence we see that  $a\|\mathbf{Z}\|_{\mathcal{M}^\perp} \leq \|\mathbf{H}\|_{\mathcal{M}^\perp} \leq \|\mathbf{Z}\|_{\mathcal{M}^\perp}$ .

## Geometric characterization of the effect of activation functions

- $\sigma$ : center  $\mathbf{Z}_{\mathcal{M}^\perp}/2$ , radius  $r := \left( \|\mathbf{Z}_{\mathcal{M}^\perp}/2\|_F^2 - \langle \mathbf{Z}_{\mathcal{M}}^+, \mathbf{Z}_{\mathcal{M}}^- \rangle_F \right)^{1/2}$ .
- $\sigma_a$ : center  $(1+a)\mathbf{Z}_{\mathcal{M}^\perp}/2$ , radius  $r_a := \left( \|(1-a)\mathbf{Z}_{\mathcal{M}^\perp}/2\|_F^2 - (1-a)^2 \langle \mathbf{Z}_{\mathcal{M}}^+, \mathbf{Z}_{\mathcal{M}}^- \rangle_F \right)^{1/2}$ .
- Prop. 1 and 2 imply the precise location of  $\mathbf{H}_{\mathcal{M}^\perp}$  (or the smoothness  $\|\mathbf{H}_{\mathcal{M}^\perp}\|_F = \|\mathbf{H}\|_{\mathcal{M}^\perp}$ ) depends on the center and the radius of the spheres. Given a fixed  $\mathbf{Z}_{\mathcal{M}^\perp}$ , the center of the spheres remains unchanged and their radii  $r$  and  $r_a$  are only affected by changes in  $\mathbf{Z}_{\mathcal{M}}$ .
- Next, we focus on analyzing how **changes in  $\mathbf{Z}_{\mathcal{M}}$  impact  $\|\mathbf{H}\|_{\mathcal{M}^\perp}$** , i.e. the smoothness of node features.

How changes in  $Z_{\mathcal{M}}$  impact  $\|H\|_{\mathcal{M}^\perp}$ ?

## Distance to the eigenspace $\mathcal{M}$

- Prop 1 and 2 show that both ReLU and leaky ReLU reduce the distance of node features to the eigenspace  $\mathcal{M}$ , i.e.  $\|\mathbf{H}\|_{\mathcal{M}^\perp} \leq \|\mathbf{Z}\|_{\mathcal{M}^\perp}$ .
- Consider  $\mathbf{Z}, \mathbf{Z}' \in \mathbb{R}^{d \times n}$  s.t.  $\mathbf{Z}_{\mathcal{M}^\perp} = \mathbf{Z}'_{\mathcal{M}^\perp}$  but  $\mathbf{Z}_{\mathcal{M}} \neq \mathbf{Z}'_{\mathcal{M}}$ . Let  $\mathbf{H}, \mathbf{H}'$  be the output of  $\mathbf{Z}, \mathbf{Z}'$  via ReLU or leaky ReLU, respectively.
  - > We have  $\|\mathbf{H}\|_{\mathcal{M}^\perp} \leq \|\mathbf{Z}\|_{\mathcal{M}^\perp}$  and  $\|\mathbf{H}'\|_{\mathcal{M}^\perp} \leq \|\mathbf{Z}'\|_{\mathcal{M}^\perp}$ .
  - >  $\mathbf{Z}_{\mathcal{M}^\perp} = \mathbf{Z}'_{\mathcal{M}^\perp}$  implies that  $\|\mathbf{Z}\|_{\mathcal{M}^\perp} = \|\mathbf{Z}'\|_{\mathcal{M}^\perp} \Rightarrow \|\mathbf{H}'\|_{\mathcal{M}^\perp} \leq \|\mathbf{Z}\|_{\mathcal{M}^\perp}$ .
- In other words, when  $\mathbf{Z}_{\mathcal{M}^\perp} = \mathbf{Z}'_{\mathcal{M}^\perp}$  is fixed, changing  $\mathbf{Z}_{\mathcal{M}}$  to  $\mathbf{Z}'_{\mathcal{M}}$  can not affect the fact that ReLU and leaky ReLU smooth node features. — Resonating with existing theories (Oono & Suzuki, ICLR 2019, Cai & Wang, arXiv:2006.13318).

## Altering the eigenspace projection

- Let  $\mathbf{z}$  be a vector with  $z_i$  being the feature of the  $i^{\text{th}}$  node, we consider

$$\mathbf{z}(\alpha) = \mathbf{z} - \alpha \mathbf{e},$$

where  $\mathbf{e}$  is the only eigenvector of  $\mathbf{G}$  associated with the eigenvalue 1.

- It is clear that

$$\mathbf{z}(\alpha)_{\mathcal{M}^\perp} = \mathbf{z}_{\mathcal{M}^\perp} \text{ and } \mathbf{z}(\alpha)_{\mathcal{M}} = \mathbf{z}_{\mathcal{M}} - \alpha \mathbf{e},$$

where we see that  $\alpha$  only alters  $\mathbf{z}_{\mathcal{M}}$  while preserves  $\mathbf{z}_{\mathcal{M}^\perp}$ .

- Consider a connected graph with 100 nodes with each being assigned a random degree between 2 to 10. Then we assign an initial node feature  $\mathbf{x} \in \mathbb{R}^{100}$ , sampled uniformly on the interval  $[-1.5, 1.5]$ , with each node feature being a scalar; we study the smoothness of node features  $\mathbf{z}_\alpha = \mathbf{x} + \alpha \mathbf{e}$ , where  $\alpha \in [-1.5, 1.5]$  is the smoothness control parameter.

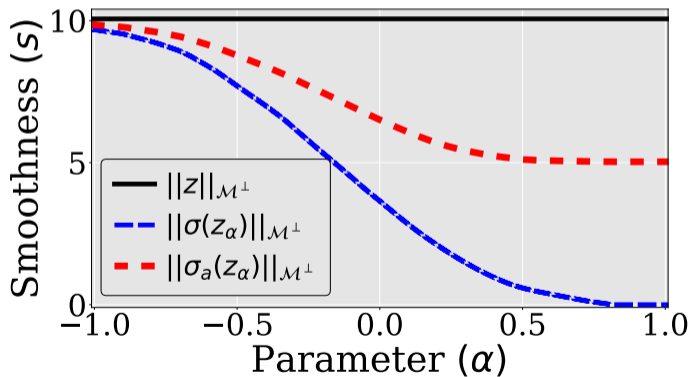


Figure: Effects of varying parameter  $\alpha$  on the smoothness of output features  $\sigma(z_\alpha)$  and  $\sigma_a(z_\alpha)$ .

## Normalized Smoothness

## Dimension-wise normalized smoothness

- For the sake of simplicity, we assume the graph is connected, i.e.  $m = 1$ .
- **Definition.** Let  $\mathbf{Z} \in \mathbb{R}^{d \times n}$  be the features over  $n$  nodes with  $\mathbf{z}^{(i)} \in \mathbb{R}^n$  ( $i = 1, \dots, d$ ) being the  $i^{\text{th}}$  row vector of  $\mathbf{Z}$ , i.e. the  $i^{\text{th}}$  dimension of the features over all nodes. Then we define the **normalized smoothness** of  $\mathbf{z}^{(i)}$  as follows:

$$s(\mathbf{z}^{(i)}) := \frac{\|\mathbf{z}_{\mathcal{M}}^{(i)}\|}{\|\mathbf{z}^{(i)}\|} \in [0, 1],$$

where we set  $s(\mathbf{z}^{(i)}) = 1$  when  $\mathbf{z}^{(i)} = \mathbf{0}$ .



## Altering the eigenspace projection

- Let  $\mathbf{z}$  be a vector with  $z_i$  being the feature of the  $i^{\text{th}}$  node, we consider

$$\mathbf{z}(\alpha) = \mathbf{z} - \alpha \mathbf{e},$$

where  $\mathbf{e}$  is the only eigenvector of  $\mathbf{G}$  associated with the eigenvalue 1.

- It is clear that

$$\mathbf{z}(\alpha)_{\mathcal{M}^\perp} = \mathbf{z}_{\mathcal{M}^\perp} \text{ and } \mathbf{z}(\alpha)_{\mathcal{M}} = \mathbf{z}_{\mathcal{M}} - \alpha \mathbf{e},$$

where we see that  $\alpha$  only alters  $\mathbf{z}_{\mathcal{M}}$  while preserves  $\mathbf{z}_{\mathcal{M}^\perp}$ .

- Consider a connected graph with 100 nodes with each being assigned a random degree between 2 to 10. Then we assign an initial node feature  $\mathbf{x} \in \mathbb{R}^{100}$ , sampled uniformly on the interval  $[-1.5, 1.5]$ , with each node feature being a scalar; we study the smoothness of node features  $\mathbf{z}_\alpha = \mathbf{x} + \alpha \mathbf{e}$ , where  $\alpha \in [-1.5, 1.5]$  is the smoothness control parameter.

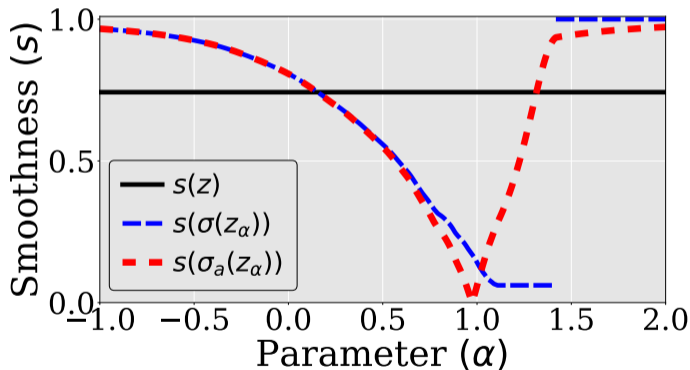
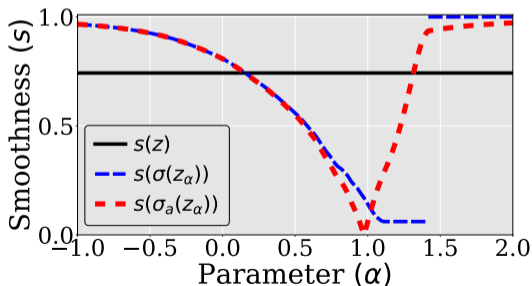


Figure: Effects of varying  $\alpha$  on the normalized smoothness of output features  $\sigma(z_\alpha)$  and  $\sigma_a(z_\alpha)$ .

**Proposition 3.** (ReLU) Suppose  $\mathbf{z}_{\mathcal{M}^\perp} \neq 0$ . Let  $\mathbf{h}(\alpha) = \sigma(\mathbf{z}(\alpha))$  with  $\sigma$  being ReLU, then

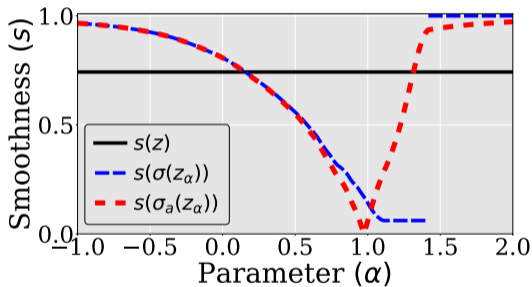
$$\min_{\alpha} s(\mathbf{h}(\alpha)) = \sqrt{\frac{\sum_{x_i = \max \mathbf{x}} d_i}{\sum_{j=1}^n d_j}} \quad \text{and} \quad \max_{\alpha} s(\mathbf{h}(\alpha)) = 1,$$

where  $\mathbf{x} := \tilde{\mathbf{D}}^{-\frac{1}{2}} \mathbf{z}$ ,  $\max \mathbf{x} = \max_{1 \leq i \leq n} x_i$ , and  $\tilde{\mathbf{D}} = \text{diag}(d_1, d_2, \dots, d_n)$ . Also, the normalized smoothness  $s(\mathbf{h}(\alpha))$  is monotone increasing as  $\alpha$  decreases whenever  $\alpha < \|\tilde{\mathbf{D}}^{\frac{1}{2}} \mathbf{u}_n\| \max \mathbf{x}$  and it has range  $[\min_{\alpha} s(\mathbf{h}(\alpha)), 1]$ .



**Figure:** Effects of varying  $\alpha$  on the normalized smoothness of output features  $\sigma(\mathbf{z}_\alpha)$  and  $\sigma_a(\mathbf{z}_\alpha)$ .

**Proposition 4.** (Leaky ReLU) Suppose  $\mathbf{z}_{\mathcal{M}^\perp} \neq 0$ . Let  $\mathbf{h}(\alpha) = \sigma_a(\mathbf{z}(\alpha))$  with  $\sigma_a$  being leaky ReLU, then 1)  $\min_\alpha s(\mathbf{h}(\alpha)) = 0$ , and 2)  $\sup_\alpha s(\mathbf{h}(\alpha)) = 1$ . Also,  $s(\mathbf{h}(\alpha))$  has range  $[0, 1)$ .



**Figure:** Effects of varying  $\alpha$  on the normalized smoothness of output features  $\sigma(\mathbf{z}_\alpha)$  and  $\sigma_a(\mathbf{z}_\alpha)$ .

**Theorem 1.** Suppose  $\mathbf{z}_{\mathcal{M}^\perp} \neq 0$ . Let  $\mathbf{h}(\alpha) = \sigma(\mathbf{z}(\alpha))$  or  $\sigma_a(\mathbf{z}(\alpha))$  with  $\sigma$  being ReLU and  $\sigma_a$  being leaky ReLU. Then we have  $\|\mathbf{z}\|_{\mathcal{M}^\perp} \geq \|\mathbf{h}(\alpha)\|_{\mathcal{M}^\perp}$  for any  $\alpha \in \mathbb{R}$ . However,  $s(\mathbf{h}(\alpha))$  can be smaller than, larger than, or equal to  $s(\mathbf{z})$  for different values of  $\alpha$ .

# Controlling the Smoothness of Node Features

## Controlling the smoothness of node features

- Our proposed smoothness control term (SCT):

$$\mathbf{B}_{\alpha^l} = \sum_{i=1}^m \alpha_i^l \mathbf{e}_i^\top,$$

where  $l$  is the layer index,  $\{\mathbf{e}_i\}_{i=1}^m$  is the orthonormal basis of the eigenspace  $\mathcal{M}$ , and  $\alpha^l$  is a collection of learnable vectors  $\{\alpha_i^l\}_{i=1}^m$  with  $\alpha_i^l \in \mathbb{R}^d$  being approximated by an MLP.

- GCN-SCT:

$$\mathbf{H}^l = \sigma(\mathbf{W}^l \mathbf{H}^{l-1} \mathbf{G} + \mathbf{B}_{\alpha^l}).$$

- GCNII-SCT:

$$\mathbf{H}^l = \sigma(((1 - \alpha_l) \mathbf{H}^{l-1} \mathbf{G} + \alpha_l \mathbf{H}^0) ((1 - \beta_l) \mathbf{I} + \beta_l \mathbf{W}^l) + \mathbf{B}_{\alpha^l}),$$

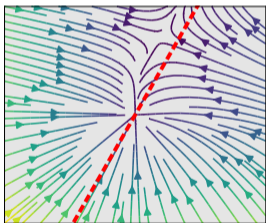
where the residual connection and identity mapping are consistent with GCNII.

## Node feature trajectory

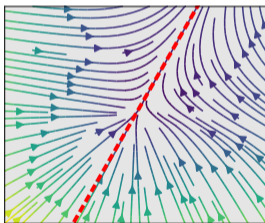
- Consider a connected graph with two nodes with 1D node features. GCL becomes

$$\mathbf{h}^1 = \sigma(w\mathbf{h}^0\mathbf{G} + \mathbf{b}_\alpha),$$

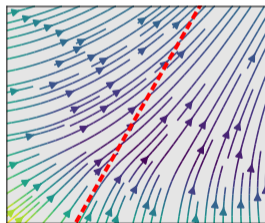
where  $w = 1.2$ ,  $\mathbf{h}^0, \mathbf{h}^1, \mathbf{b}_\alpha \in \mathbb{R}^2$ , and  $\mathbf{G} \in \mathbb{R}^{2 \times 2}$ . We select a positive definite matrix  $\mathbf{G}$  with the largest eigenvalue 1;  $\mathbf{G}$  is defined to be  $[0.592, 0.194; 0.194, 0.908]$ . Twenty initial node feature vectors  $\mathbf{h}^0$  are sampled evenly in the domain  $[-1, 1] \times [-1, 1]$ .



a)  $\alpha = -0.25$



b)  $\alpha = 0.0$



c)  $\alpha = 1.0$

**Figure:** Node feature trajectories, with colored magnitude, for varying smoothness control parameter  $\alpha$ . For classical GCN b), the node features converge to the eigenspace  $\mathcal{M}$  (red dashed line).



Layers	2	4	16	32
<b>Cora</b>				
GCN/GCN-SCT	81.1/ <b>82.9</b>	80.4/ <b>82.8</b>	64.9/ <b>71.4</b>	60.3/ <b>67.2</b>
GCNII/GCNII-SCT	82.2/ <b>83.8</b>	82.6/ <b>84.3</b>	84.6/ <b>84.8</b>	85.4/ <b>85.5</b>
EGNN/EGNN-SCT	83.2/ <b>84.1</b>	84.2/ <b>84.5</b>	<b>85.4</b> /83.3	<b>85.3</b> /82.0
<b>Citeseer</b>				
GCN/GCN-SCT	<b>70.3</b> /69.9	67.6/ <b>67.7</b>	18.3/ <b>55.4</b>	25.0/ <b>51.0</b>
GCNII/GCNII-SCT	68.2/ <b>72.8</b>	68.9/ <b>72.8</b>	72.9/ <b>73.8</b>	<b>73.4</b> / <b>73.4</b>
EGNN/EGNN-SCT	72.0/ <b>73.1</b>	71.9/ <b>72.0</b>	72.4/ <b>72.6</b>	72.3/ <b>72.9</b>
<b>PubMed</b>				
GCN/GCN-SCT	79.0/ <b>79.8</b>	76.5/ <b>78.4</b>	40.9/ <b>76.1</b>	22.4/ <b>77.0</b>
GCNII/GCNII-SCT	78.2/ <b>79.7</b>	78.8/ <b>80.1</b>	80.2/ <b>80.7</b>	79.8/ <b>80.7</b>
EGNN/EGNN-SCT	79.2/ <b>79.8</b>	79.5/ <b>80.4</b>	80.1/ <b>80.3</b>	80.0/ <b>80.4</b>
<b>Coauthor-Physics</b>				
GCN/GCN-SCT	92.4/ <b>92.6</b>	92.1/ <b>92.5</b>	13.5/ <b>50.9</b>	13.1/ <b>43.6</b>
GCNII/GCNII-SCT	92.5/ <b>94.4</b>	92.9/ <b>94.2</b>	92.9/ <b>93.7</b>	92.9/ <b>94.1</b>
EGNN/EGNN-SCT	92.6/ <b>93.9</b>	92.9/ <b>94.1</b>	93.1/ <b>94.0</b>	93.3/ <b>93.8</b>
<b>Ogbn-arxiv</b>				
GCN/GCN-SCT	70.4/ <b>72.1</b>	71.7/ <b>72.7</b>	70.6/ <b>72.3</b>	68.5/ <b>72.3</b>
GCNII/GCNII-SCT	70.1/ <b>72.0</b>	71.4/ <b>72.1</b>	71.5/ <b>72.4</b>	70.5/ <b>72.1</b>
EGNN/EGNN-SCT	68.4/ <b>68.5</b>	71.1/ <b>71.3</b>	72.7/ <b>72.8</b>	<b>72.7</b> /72.3

**Table:** Test accuracy for models of varying depth on citation networks with benchmark splits. (Unit:%)

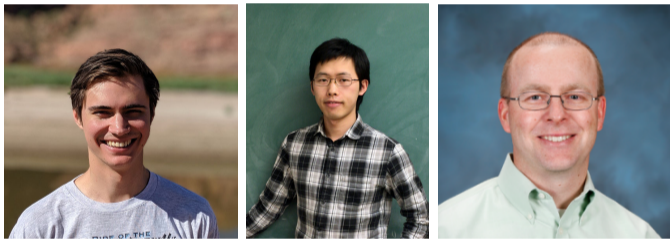
Cornell	Texas	Wisconsin	Chameleon
52.70/ <b>55.95</b> (0.007/0.018)	52.16/ <b>62.16</b> (0.007/0.008)	45.88/ <b>54.71</b> (0.007/0.008)	28.18/ <b>38.44</b> (0.006/0.007)
74.86/ <b>75.41</b> (0.020/0.020)	69.46/ <b>83.34</b> (0.031/0.020)	74.12/ <b>86.08</b> (0.020/0.015)	60.61/ <b>64.52</b> (0.015/0.013)

**Table:** Mean test accuracy results and average computational time per epoch (in the parenthesis) for the WebKB and WikipediaNetwork datasets with fixed 48/32/20% splits. First row: GCN/GCN-SCT. Second row: GCNII/GCNII-SCT. (Unit:% (second))

## References

Shih-Hsin Wang\*, Justin Baker\*, Cory Hauck, and Bao Wang, The Effects of Activation Functions on the Over-smoothing Issue of Graph Convolutional Networks, preprint, 2023.

# Implicit Graph Neural Networks: A Monotone Operator Viewpoint



---

Justin Baker\*, Qingsong Wang\*, Cory Hauck, and Bao Wang, Implicit Graph Neural Networks: A Monotone Operator Viewpoint, ICML, 2023.

## Implicit GNNs

- Implicit GNN (IGNN)

$$\mathbf{Z}^{(k+1)} = \sigma(\mathbf{WZ}^{(k)}\mathbf{G} + g_{\mathbf{B}}(\mathbf{X})), \text{ for } k = 0, 1, 2, \dots,$$

where  $g_{\mathbf{B}}$  is a function parameterized by  $\mathbf{B}$ , e.g.  $g_{\mathbf{B}}(\mathbf{X}) = \mathbf{BXG}$  with  $\mathbf{B} \in \mathbb{R}^{d \times d}$ .

- Finding the fixed point  $\mathbf{Z}^*$  as the representation of input graph.

## Issue 1: Well-posedness of IGNN limits its expressivity of IGNN

- Well-posedness, i.e. the fixed point exists and is unique

$$\lambda_1(|\mathbf{W}|) < 1.$$

Or all eigenvalues of  $\mathbf{W}$  are less than one in magnitude.

- The selection of  $\mathbf{W}$  is limited, limiting the expressivity of IGNN.

---

Notice that all eigenvalues of  $\mathbf{G} = \hat{\mathbf{A}}$  are in  $[-1, 1]$  with  $\lambda_1(\mathbf{G}) = 1$ .

## Issue 2: When IGNN learns long-range dependencies (LRD)

- Learning LRD: each node can aggregate information from the nodes that are far apart.
- To learn LRD,  $\lambda_1(|\mathbf{W}|)$  needs to be close to one in magnitude; otherwise, the Picard iteration converges too fast, and each node only aggregates nearby nodes' features.
- Training IGNN with  $\lambda_1(|\mathbf{W}|) \rightarrow 1$ , starting from random initialization, may not happen.
- Picard iteration converges slowly when  $\lambda_1(|\mathbf{W}|) \rightarrow 1$

## A monotone operator theory viewpoint of IGNN

- Notice that  $\mathbf{Z}^{(k+1)} = \sigma(\mathbf{WZ}^{(k)}\mathbf{G} + g_{\mathbf{B}}(\mathbf{X}))$  can be rewritten as the following vectorized equation

$$\text{vec}(\mathbf{Z}^{(k+1)}) = \sigma(\mathbf{G}^{\top} \otimes \mathbf{W} \text{vec}(\mathbf{Z}^{(k)}) + \text{vec}(g_{\mathbf{B}}(\mathbf{X}))), \quad (1)$$

where  $\mathbf{G}^{\top} \otimes \mathbf{W}$  denotes the Kronecker product between  $\mathbf{G}$  and  $\mathbf{W}$ .



## A monotone operator theory viewpoint of IGNN

- Finding a fixed point of (1) is equivalent to solving the monotone inclusion problem

$$\text{find } 0 \in (\mathcal{F} + \mathcal{G})(\text{vec}(\mathbf{Z})^*),$$

where

$$\mathcal{F}(\text{vec}(\mathbf{Z})) = (\mathbf{I} - \mathbf{G}^\top \otimes \mathbf{W})\text{vec}(\mathbf{Z}) - \text{vec}(\mathbf{g}_B(\mathbf{X})) \quad \text{and} \quad \mathcal{G} = \partial f,$$

where  $f$  is a convex closed proper (CCP) function such that

$$\sigma(x) = \text{prox}_f^1(x) = \underset{z}{\text{argmin}} \left\{ \frac{1}{2} \|x - z\|^2 + f(z) \right\}.$$

- Notice that when  $\sigma$  is ReLU, then  $\sigma = \text{prox}_f^\alpha$  for  $\forall \alpha > 0$  with  $f$  being the indicator of the positive octant, i.e.  $f(x) = I\{x \geq 0\}$ .

## Well-posedness of MIGNN

- MIGNN: monotone operator theory viewpoint of IGNN.
- The fixed point  $\mathbf{Z}^*$  exists and is unique if  $\mathcal{F}$  is strongly monotone.
- If  $\mathbf{I} - \mathbf{G}^\top \otimes \mathbf{W} \succeq m\mathbf{I}$  for some  $m > 0$ , then  $\mathcal{F}$  is strongly monotone.

## Monotone parameterization of MIGNN: Enhancing expressivity of IGNN

- We consider the following MIGNN model

$$\mathbf{Z}^{(k+1)} = \sigma(\mathbf{W}\mathbf{Z}^{(k)}\mathbf{G} + g_B(\mathbf{X})).$$

- We let  $\mathbf{G} = \frac{\mathbf{L}}{2}$  where  $\mathbf{L} := \mathbf{D}^{-1/2}(\mathbf{D} - \mathbf{A})\mathbf{D}^{-1/2}$  is the normalized Laplacian.

- We parameterize  $\mathbf{W}$  with the following [monotone parameterization](#)

$$\mathbf{W} = (1 - m)\mathbf{I} - \mathbf{C}\mathbf{C}^\top + \mathbf{F} - \mathbf{F}^\top,$$

where  $\mathbf{C}, \mathbf{F} \in \mathbb{R}^{d \times d}$  are arbitrary matrices, and  $m > 0 \in \mathbb{R}$ .

## Monotone parameterization of MIGNN: Enhancing expressivity of IGNN

- The monotone parameterization guarantees the operator  $\mathcal{F}$  to be strongly monotone.
- The monotone parameterization allows the eigenvalues of  $\mathbf{W}$  to be much less than  $-1$ , which is more flexible than IGNN.

## Orthogonal parameterization of MIGNN: Stabilizing learning LRD

- Consider the following MIGNN model

$$\mathbf{Z}^{(k+1)} = \sigma(\mathbf{WZ}^{(k)}\mathbf{G} + g_B(\mathbf{X})).$$

- We parameterize  $\mathbf{W}$  using the following scaled Cayley map

$$\mathbf{W} = \phi(\gamma)(\mathbf{I} - \mathbf{S})(\mathbf{I} + \mathbf{S})^{-1},$$

where  $\phi(\cdot)$  is the sigmoid function.  $\mathbf{S} = \mathbf{C} - \mathbf{C}^\top$  is a skew-symmetric matrix with  $\mathbf{C} \in \mathbb{R}^{d \times d}$  an arbitrary matrix.

- Notice that the matrix  $(\mathbf{I} - \mathbf{S})(\mathbf{I} + \mathbf{S})^{-1}$  is orthogonal.

## Finding the fixed point of the equilibrium equation

- Picard iteration may not converge for MIGNN with monotone parameterization, i.e.,  $\mathbf{W} = (1 - m)\mathbf{I} - \mathbf{C}\mathbf{C}^\top + \mathbf{F} = \mathbf{F}^\top$ .
- Picard iteration suffers from slow convergence for MIGNN with orthogonal parameterization, i.e.,  $\mathbf{W} = (\mathbf{I} - \mathbf{S})(\mathbf{I} + \mathbf{S})^{-1}$  with  $\mathbf{S} = \mathbf{C} - \mathbf{C}^\top$ .
- Need new algorithms to find the fixed point of MIGNN.

## Forward-backward splitting (FB): MIGNN with monotone parameterization

- Finding the fixed point of MIGNN,  $\mathbf{Z}^{(k+1)} = \sigma(\mathbf{WZ}^{(k)}\mathbf{G} + g_B(\mathbf{X}))$ , with monotone parameterization

$$\mathbf{Z}^{(k+1)} := F_\alpha^{\text{FB}}(\mathbf{Z}^{(k)}) := \text{prox}_f^\alpha \left( \mathbf{Z}^{(k)} - \alpha \cdot \left( \mathbf{Z}^{(k)} - \mathbf{WZ}^{(k)}\mathbf{G} - g_B(\mathbf{X}) \right) \right),$$

where  $\alpha > 0$  is an appropriate constant.

- $$\begin{aligned} \mathbf{Z}^{(k+1/2)} &= \mathbf{Z}^{(k)} - \alpha \cdot \left( \mathbf{Z}^{(k)} - \mathbf{WZ}^{(k)}\mathbf{G} - g_B(\mathbf{X}) \right) \\ \mathbf{Z}^{(k+1)} &= \text{prox}_f^\alpha(\mathbf{Z}^{(k+1/2)}). \end{aligned}$$

- Resulting the model MIGNN-Mon.

## Peaceman-Rachford splitting (PR): MIGNN with orthogonal parameterization

- PR finds the solution  $\mathbf{Z}^*$  of the MIGNN by letting

$$\mathbf{Z}^* = \text{prox}_f^\alpha(\mathbf{U}^*),$$

where  $\mathbf{U}^*$  is the solution of the following fixed point iterations:

$$\text{vec}(\mathbf{U}^{(k+1)}) = F_\alpha^{\text{PR}}(\text{vec}(\mathbf{U}^{(k)})) := \mathcal{C}_\mathcal{F}\mathcal{C}_\mathcal{G}(\text{vec}(\mathbf{U}^{(k)})),$$

where

$$\mathcal{R}_\mathcal{T} = (\mathcal{I} + \alpha\mathcal{T})^{-1},$$

and

$$\mathcal{C}_\mathcal{T} = 2\mathcal{R}_\mathcal{T} - \mathcal{I}.$$



## Peaceman-Rachford splitting (PR): MIGNN with orthogonal parameterization

- Let  $\mathbf{u}^k := \text{vec}(\mathbf{U}^{(k)})$ , then we can formulate PR as follows

$$\mathbf{u}^{k+1} := F_{\alpha}^{\text{PR}}(\mathbf{u}^k) = 2\mathbf{V} \left( 2\text{prox}_f^{\alpha}(\mathbf{u}^k) - \mathbf{u}^k + \alpha \text{vec}(\mathbf{g}_{\mathbf{B}}(\mathbf{X})) \right) - 2\text{prox}_f^{\alpha}(\mathbf{u}^k) + \mathbf{u}^k,$$

where the matrix  $\mathbf{V} := (\mathbf{I} + \alpha(\mathbf{I} - \mathbf{G}^{\top} \otimes \mathbf{W}))^{-1}$  and  $\mathbf{u}^0$  is the zero vector.

- Computing  $\mathbf{V}(\mathbf{x}^k)$  is expensive:
  - > Bartels–Stewart algorithm, which requires diagonalizing the matrix  $\mathbf{G}$  and  $\mathbf{W}$ .

## PR with Neumann series approximation

- Notice that

$$\begin{aligned}\mathbf{V}(\mathbf{u}^k) &= (\mathbf{I} + \alpha(\mathbf{I} - \mathbf{G}^\top \otimes \mathbf{W}))^{-1}(\mathbf{u}^k) \\ &= \frac{1}{1 + \alpha} \left( \mathbf{I} - \frac{\mathbf{G}^\top \otimes \mathbf{W}}{1 + 1/\alpha} \right)^{-1} (\mathbf{u}^k) \\ &= \frac{1}{1 + \alpha} \sum_{i=0}^{\infty} \frac{\text{vec}(\mathbf{W}^i \mathbf{U}^{(k)} \mathbf{G}^i)}{(1 + 1/\alpha)^i}.\end{aligned}$$

- $K$ -th order Neumann series approximation of  $\mathbf{V}(\mathbf{u}^k)$ :

$$\mathbf{N}_K(\text{vec}(\mathbf{U}^k)) := \frac{1}{1 + \alpha} \sum_{i=0}^K \frac{\text{vec}(\mathbf{W}^i \mathbf{U}^k \mathbf{G}^i)}{(1 + 1/\alpha)^i}.$$

- $K$ -th order Neumann series approximation of PR iteration

$$\mathbf{u}^{k+1} := \tilde{F}_\alpha^{\text{PR},K}(\mathbf{u}^k) = 2\mathbf{N}_K\left(2\text{prox}_f^\alpha(\mathbf{u}^k) - \mathbf{u}^k + \alpha \text{vec}(g_{\mathbf{B}}(\mathbf{X}))\right) - 2\text{prox}_f^\alpha(\mathbf{u}^k) + \mathbf{u}^k.$$

## MIGNN with diffusion convolution

- We can set  $\mathbf{G}$  to be the combination of higher powers of  $\hat{\mathbf{A}}$  or  $\mathbf{L}$ , making each node to aggregate multi-hops neighbors' features in each iteration.
- We let  $\mathbf{G} = \tilde{\mathbf{D}}^{-1/2}(\mathbf{A} + \dots + \mathbf{A}^P)\tilde{\mathbf{D}}^{-1/2}$  for any positive integer  $P$ , where  $\tilde{\mathbf{D}}$  is the degree matrix with  $\tilde{D}_{ii} = \sum_{j=1}^n \sum_{k=1}^P (\mathbf{A}^k)_{ij}$ .
- MIGNN with  $P$ -th order diffusion matrix  $\mathbf{G}$

$$\mathbf{Z} = \sigma(\mathbf{WZ}\tilde{\mathbf{D}}^{-1/2}(\mathbf{A} + \mathbf{A}^2 + \dots + \mathbf{A}^P)\tilde{\mathbf{D}}^{-1/2} + g_{\mathbf{B}}(\mathbf{X})).$$

- We denote the model as **MIGNN-NKDP** when it is implemented by using the  $P$ -th order diffusion and the  $K$ -th order Neumann series approximated PR iteration.

## Directed chain classification

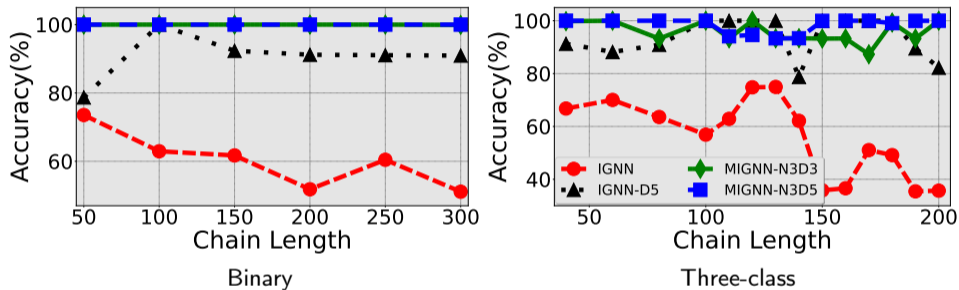


Figure: The accuracy of IGNN and MIGNN for classifying directed chains of different lengths.

## Directed chain classification

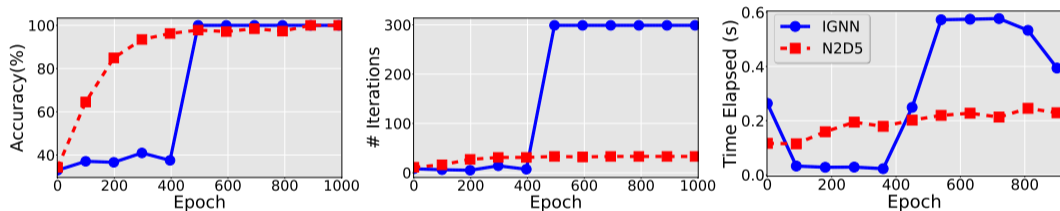


Figure: MIGNN-N2D5 vs. IGNN for three class chains classification (length: 140).

## Graph node classification: Citation networks

Datasets	Cora	Citeseer	Pubmed
Geom-GCN	85.27	<b>77.99</b>	<b>90.05</b>
GCNII	<b>88.49</b>	77.08	89.57
APPNP	85.09	75.73	79.73
GCN+GDC	83.58	73.35	78.72
GIND	88.25	76.81	89.22
IGNN	85.80	75.24	87.66
EIGNN (Ours)	85.89	75.31	87.92
MIGNN-Mon (Ours)	86.82	76.59	88.00
MIGNN-N5D1	87.04	74.91	83.55

**Table:** Node classification mean accuracy (%) for 10-fold cross-validation.

## Graph classification: bioinformatics-related tasks

Datasets	MUTAG	PTC	COX2	PROTEINS	NCI1
# graphs/Avg # nodes	188/17.9	344/25.5	467/41.2	1113/39.1	4110/29.8
WL	84.1 ± 1.9	58.0 ± 2.5	83.2 ± 0.2	74.7 ± 0.5	84.5 ± 0.5
DCNN	67.0	56.6	—	61.3	62.6
DGCNN	85.8	58.6	—	75.5	74.4
GIN	89.4 ± 5.6	64.6 ± 7.0	—	76.2 ± 3.4	82.7 ± 1.7
FDGNN	88.5 ± 3.8	63.4 ± 5.4	83.3 ± 2.9	76.8 ± 2.9	77.8 ± 1.6
IGNN	76.0 ± 13.4	60.5 ± 6.4	79.7 ± 3.4	76.5 ± 3.4	73.5 ± 1.9
GIND	89.3 ± 7.4	66.9 ± 6.6	84.8 ± 4.2	77.2 ± 2.9	78.8 ± 2.9
GSN	92.2 ± 7.5	68.2 ± 7.2	—	76.6 ± 5.0	83.5 ± 2.0
SIN	—	—	—	76.5 ± 3.3	82.8 ± 2.2
CIN	92.7 ± 6.1	68.2 ± 5.6	—	77.0 ± 4.3	83.6 ± 1.4
MIGNN-Mon	81.8 ± 9.1	72.6 ± 6.7	85.0 ± 5.3	77.9 ± 3.4	73.6 ± 2.0
MIGNN-N1D1	86.1 ± 9.1	70.9 ± 6.5	86.5 ± 2.8	79.0 ± 3.3	78.4 ± 1.2
MIGNN-N3D1	91.4 ± 7.5	71.2 ± 3.2	88.2 ± 4.1	80.1 ± 3.8	80.8 ± 1.81

**Table:** Graph classification mean accuracy (%) ± standard deviation for 10-fold cross-validation.

## Graph classification: bioinformatics-related tasks

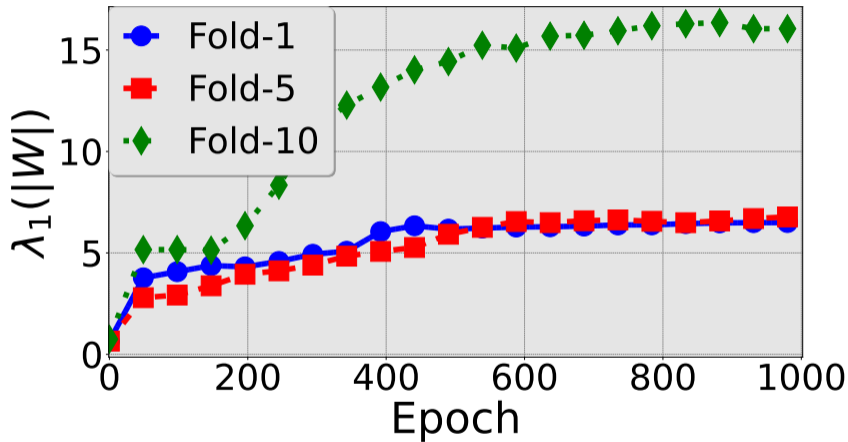


Figure:  $\lambda_1(|W|)$  of MIGNN-Mon vs. Epoch on MUTAG.



## Summary

I. How activation functions affect the smoothness of node features.

I.1 Geometric characterization

I.2 Smoothness control

II. Monotone operator-based implicit graph neural networks

I.1 Stable and accurate graph deep learning

I.2 Fast convergence and learning long-range dependencies