

Break the Ceiling: Stronger Multi-scale Deep Graph Convolutional Networks

Sitao Luan & Mingde Zhao^{*}, Xiao-Wen Chang^{*}, Doina Precup[†]

Abstract

We analyze how existing Graph Convolutional Networks (GCNs) have limited expressive power due to the constraint of the activation functions and their architectures. We generalize spectral graph convolution and deep GCN in block Krylov subspace forms and devise 2 architectures, both with the potentials to be extended deeper but each making use of the multi-scale information differently. On several node classification tasks, with or without validation set, the 2 proposed architectures achieve SOTA performance.

Problem of GCN: activation function

Deepened GCN can be described as

$\mathbf{Y}' = \mathbf{L} \text{ReLU}(\cdots \mathbf{L} \text{ReLU}(\mathbf{L} \text{ReLU}(\mathbf{L} \mathbf{X} \mathbf{W}_0) \mathbf{W}_1) \mathbf{W}_2 \cdots) \mathbf{W}_n$
where $\mathbf{L} = \tilde{\mathbf{D}}^{-1/2} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-1/2}$, $\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{I}$ and $\tilde{\mathbf{D}} = \sum_j \tilde{A}_{ij}$, $\mathbf{A} \in \mathbb{R}^{N \times N}$ is the adjacency matrix.

Theorem 1. Suppose the graph \mathcal{G} has k connected components and the diffusion operator L is defined as above. Let $\mathbf{X} \in \mathbb{R}^{N \times F}$ be any block vector and let W_j be any non-negative parameter matrix with $\|W_j\|_2 \leq 1$ for $j = 0, 1, \dots$. If \mathcal{G} has no bipartite components, then we have $\lim_{n \rightarrow \infty} \text{rank}(\mathbf{Y}') \leq k$.

Conjecture 1. Theorem 1 still holds without the non-negative constraint on the parameter matrices (see Figure 1).

Theorem 2. Suppose the n -dimensional \mathbf{x} and \mathbf{y} are independently sampled from a continuous distribution and the activation function Tanh is applied to $[\mathbf{x}, \mathbf{y}]$ pointwisely, then

$$\mathbb{P}(\text{rank}(\text{Tanh}([\mathbf{x}, \mathbf{y}])) = \text{rank}([\mathbf{x}, \mathbf{y}])) = 1$$

Problem of GCN: architecture

Definition Let \mathbb{S} be a block vector subspace of $\mathbb{R}^{F \times F}$. Given a set of block vectors $\{\mathbf{X}_k\}_{k=1}^m \subset \mathbb{R}^{N \times F}$, the \mathbb{S} -span of $\{\mathbf{X}_k\}_{k=1}^m$ is defined as $\text{span}^{\mathbb{S}}\{\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_m\} := \{\sum_{k=1}^m \mathbf{X}_k C_k : C_k \in \mathbb{S}\}$.

The order- m block Krylov subspace w.r.t. $\mathbf{A} \in \mathbb{R}^{N \times N}$, $\mathbf{B} \in \mathbb{R}^{N \times F}$ and \mathbb{S} is defined as $\mathcal{K}_m^{\mathbb{S}}(\mathbf{A}, \mathbf{B}) := \text{span}^{\mathbb{S}}\{\mathbf{B}, \mathbf{AB}, \dots, \mathbf{A}^{m-1}\mathbf{B}\}$. The corresponding block Krylov matrix is defined as $K_m(\mathbf{A}, \mathbf{B}) := [\mathbf{B}, \mathbf{AB}, \dots, \mathbf{A}^{m-1}\mathbf{B}]$.

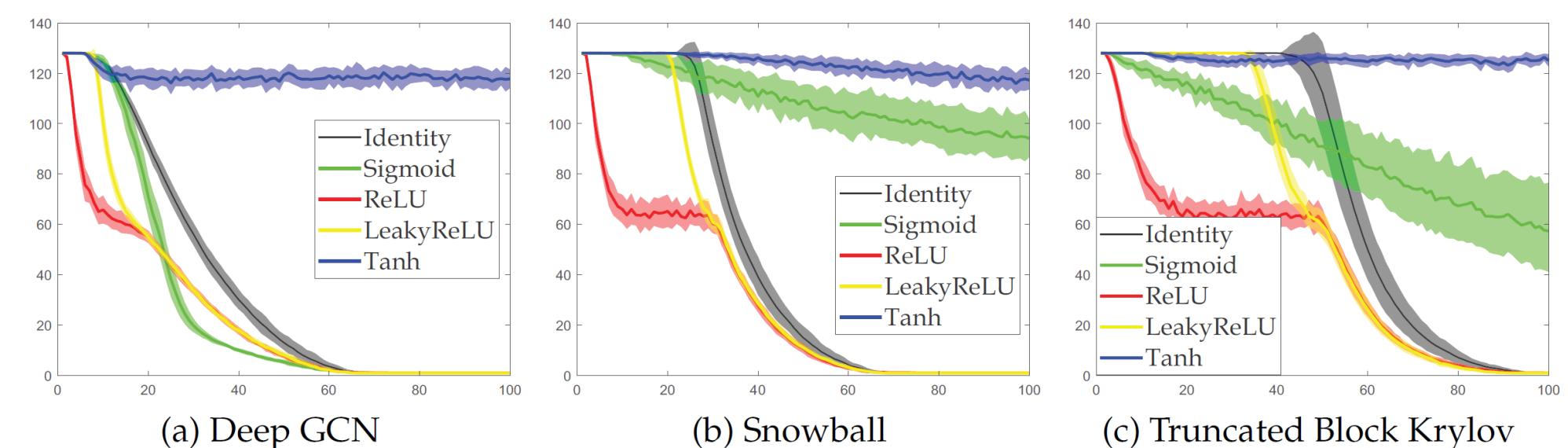


Figure 1: Changes in the number of independent features with the increment of network depth

Block Krylov Forms

Graph convolution with any well-defined analytic spectral filter g defined on L can be written as the product of a block Krylov matrix with a specific parameter matrix,

$$g(L)\mathbf{X} = \sum_{n=0}^{\infty} \frac{g^{(n)}(0)}{n!} L^n \mathbf{X} = [\mathbf{X}, L\mathbf{X}, L^2\mathbf{X}, \dots] \left[\frac{g^{(n)}(0)}{0!} I_F, \frac{g^{(n)}(0)}{1!} I_F, \frac{g^{(n)}(0)}{2!} I_F, \dots \right]^T = A'B'$$

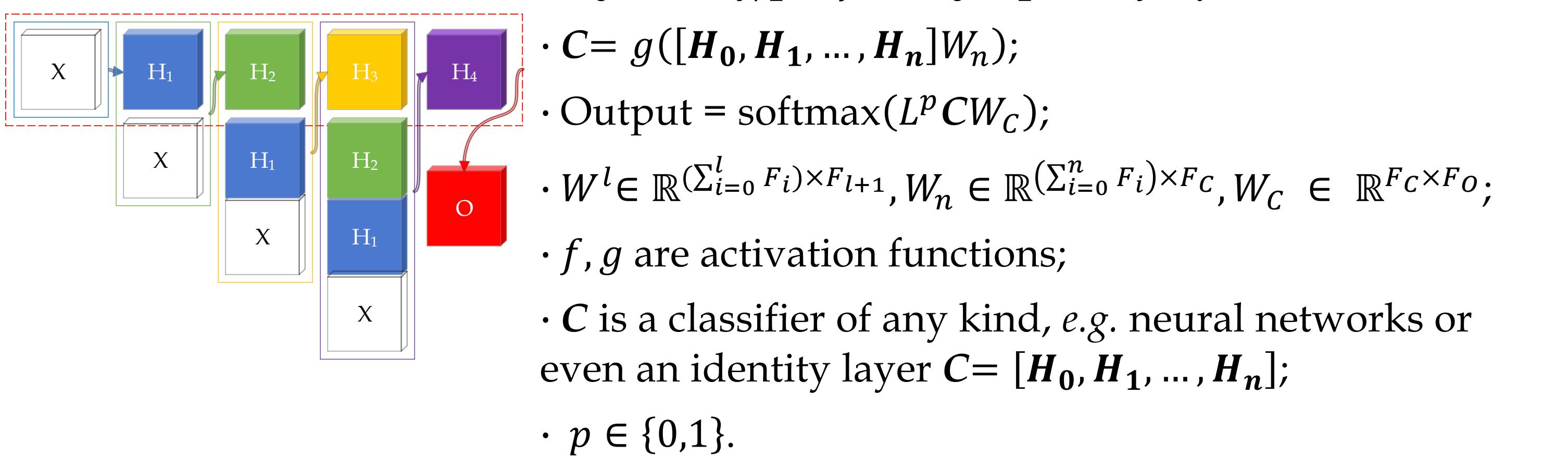
where $\rho(L) < R$, R is the convergence radius and $\rho(L)$ denotes the spectrum radius of L . $\text{Range}(A'B') = \text{Range}(A')$ and exists a smallest m s.t. $\text{span}^{\mathbb{S}}\{\mathbf{X}, L\mathbf{X}, L^2\mathbf{X}, \dots\} = \text{span}^{\mathbb{S}}\{\mathbf{X}, L\mathbf{X}, \dots, L^{m-1}\mathbf{X}\}$. Then, the convolution in hidden layer can be written as $g(L)\mathbf{X}\mathbf{W}' = K_m(L, \mathbf{X})\mathbf{W}'$. A deep graph network can be written as

$$\mathbf{Y} = \text{softmax}\{K_m(L, \mathbf{H}_n)\mathbf{W}_n^{\mathbb{S}_n}\}$$

where $\mathbf{H}_0 = \mathbf{X}, \mathbf{H}_{i+1} = h_i\{K_{m_i}(L, \mathbf{H}_i)\mathbf{W}_i^{\mathbb{S}_i}\}$, $m_i = m(L, \mathbf{H}_i)$ and h_i is activation function for $i = 0, \dots, n - 1$.

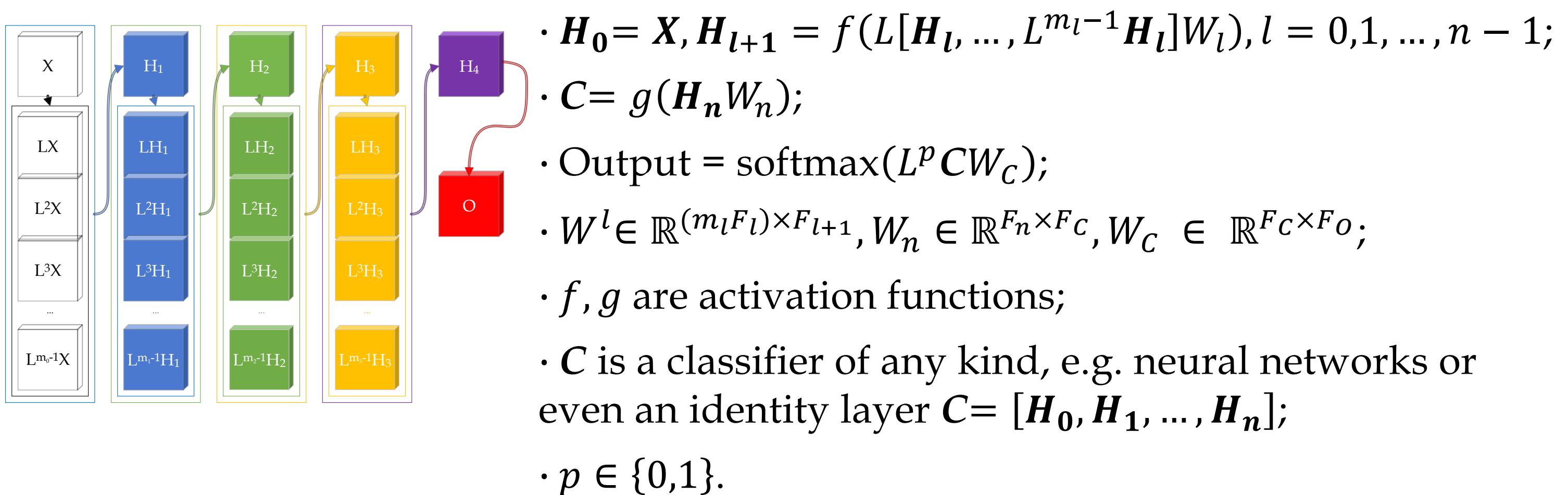
There exist some difficulties when implementing block Krylov subspace algorithm during training, but it provides an inspiration to construct a deep graph network: we need to concatenate multi-scale information in each hidden layer so that multi-hop neighborhood features can be kept for each node.

Proposed Architecture: Snowball



Snowball GCN stacks multi-scale information incrementally.

Proposed Architecture: Truncated Krylov



In truncated Krylov GCN, the local information will not be diluted because in each layer l , we start the concatenation from $L^0\mathbf{H}_l$. It can be shown that snowball GCN with linear activation function is equivalent to a one-layer truncated block Krylov network. The 2 architectures have some inherent connections for feature extractions. Truncated Krylov GCN stacks multi-scale information directly in each layer.

Experiments: Node Classification

Algorithms	0.5%	1%	3%	5.2%	0.5%	1%	3.6%	0.03%	0.05%	0.1%	0.3%
Cheby	33.9	44.2	62.1	78.0	45.3	59.4	70.1	45.3	48.2	55.2	69.8
GCN-FP	50.5	59.6	73.1	74.6	43.9	54.3	61.5	56.2	63.2	70.3	76.0
GGNN	48.2	60.5	73.1	77.6	44.3	56.0	64.6	55.8	63.3	70.4	75.8
DCNN	59.0	66.4	76.7	79.7	53.1	62.2	69.4	60.9	66.7	73.1	76.8
MPNN	46.5	56.7	72.0	78.0	41.8	54.3	64.0	53.9	59.6	67.3	75.6
GraphSAGE	37.5	49.0	64.2	74.5	33.8	51.0	67.2	45.4	53.0	65.4	76.8
GAT	41.4	48.6	56.8	83.0	38.2	46.5	72.5	50.9	50.4	59.6	79.0
GCN	50.9	62.3	76.5	80.5	43.6	55.3	68.7	57.1	64.6	73.0	77.8
LNet	58.1	66.1	77.3	79.5	53.2	61.3	66.2	60.4	68.8	73.4	78.3
AdaLN	60.8	67.5	77.7	80.4	53.8	63.3	68.7	61.0	66.0	72.8	78.1
linear Snowball	71.7	75.2	81.5	83.7	61.2	67.3	73.4	71.7	73.2	75.7	79.2
Snowball	72.9	76.2	81.1	83.4	62.4	67.0	73.2	70.9	73.0	76.1	79.5
truncated Krylov	74.1	77.6	81.8	83.5	65.3	68.2	74.2	71.5	73.2	77.0	80.1

Algorithms	0.5%	1%	2%	3%	4%	5%	0.5%	1%	2%	3%	4%	5%	0.03%	0.05%	0.1%	0.3%
LP	56.4	62.3	65.4	67.5	69.0	70.2	34.8	40.2	43.6	45.3	46.4	47.3	61.4	66.4	65.4	66.8
Cheby	38.0	52.0	62.4	70.8	74.1	77.6	31.7	42.8	59.9	66.2	68.3	69.3	40.4	47.3	51.2	72.8
Co-training	56.6	66.4	73.5	75.9	78.9	80.8	47.3	55.7	62.1	62.5	64.5	65.5	62.2	68.3	72.7	78.2
Self-training	53.7	66.1	73.8	77.2	79.4	80.0	43.3	58.1	68.2	69.8	70.4	71.0	51.9	58.7	66.8	77.0
Union	58.5	69.9	75.8	78.5	80.4	81.7	46.3	59.1	66.7	67.6	68.2	68.2	58.4	64.0	70.7	79.2
Intersection	49.7	65.0	72.9	77.1	79.4	80.2	42.9	59.1	68.8	70.1	70.8	71.2	52.0	59.3	69.7	77.6
MultiStage	61.1	63.7	74.4	76.1	77.2	77.2	53.0	57.8	63.8	68.0	69.0	70.0	57.4	64.3	70.2	77.0
M3S	61.5	67.2	75.6	77.8	78.0	78.0	56.1	62.1	66.4	70.3	70.5	70.5	59.2	64.4	70.6	76.6
GCN	42.6	56.9	67.8	74.9	77.6	79.3	33.4	46.5	62.6	66.9	68.7	69.6	46.4	49.7	56.3	76.6
GCN-SVAT	43.6	53.9	71.4	75.6												

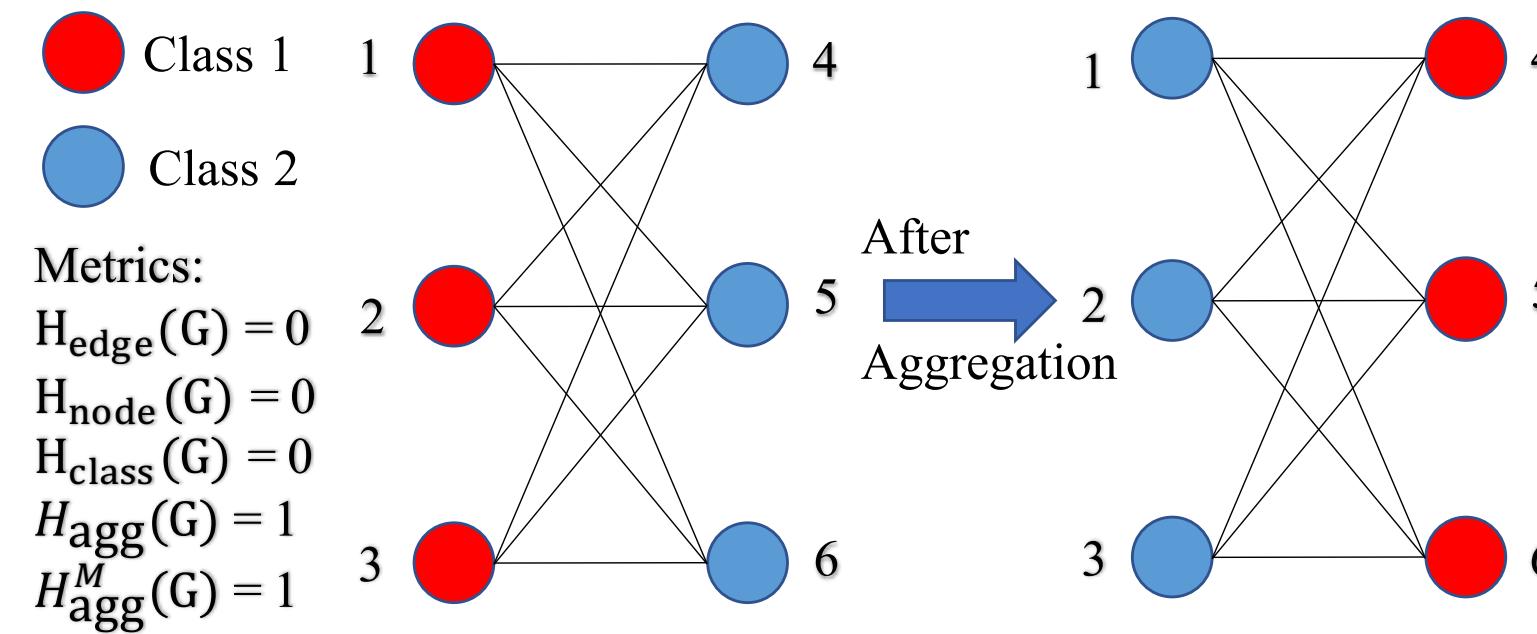
Revisiting Heterophily For Graph Neural Networks

Sitao Luan,[†] Chenqing Hua,[†] Qincheng Lu,[†] Jiaqi Zhu,[†] Mingde Zhao,[‡] Shuyuan Zhang,[‡] Xiao-Wen Chang,[‡] Doina Precup[‡]

Abstract

We first revisit the widely used homophily metrics and point out their shortcomings. Then, we study heterophily from the perspective of post-aggregation node similarity and define new homophily metrics, which are more informative. Based on this investigation, we prove that some harmful cases of heterophily can be effectively addressed by local diversification operation. Then, we propose the Adaptive Channel Mixing (ACM), a framework to adaptively exploit aggregation, diversification and identity channels node-wisely for diverse node heterophily situations. When evaluated on 10 benchmark node classification tasks, ACM-augmented baselines consistently achieve significant performance gain, exceeding state-of-the-art GNNs on most tasks without incurring significant computational burden.

Shortcomings of Homophily Metrics



Homophily: the edges tend to connect nodes from the same class, and it only describes graph-label consistency. The inconsistency relation, i.e., low homophily or heterophily, is implied to have a negative effect to the performance of GNNs, because nodes from different classes will be falsely mixed and become indistinguishable. Nevertheless, it is not always the case, e.g., the bipartite graph shown above is highly heterophilic after mean aggregation, the nodes in classes 1 and 2 just exchange colors and are still distinguishable.

Aggregation Homophily

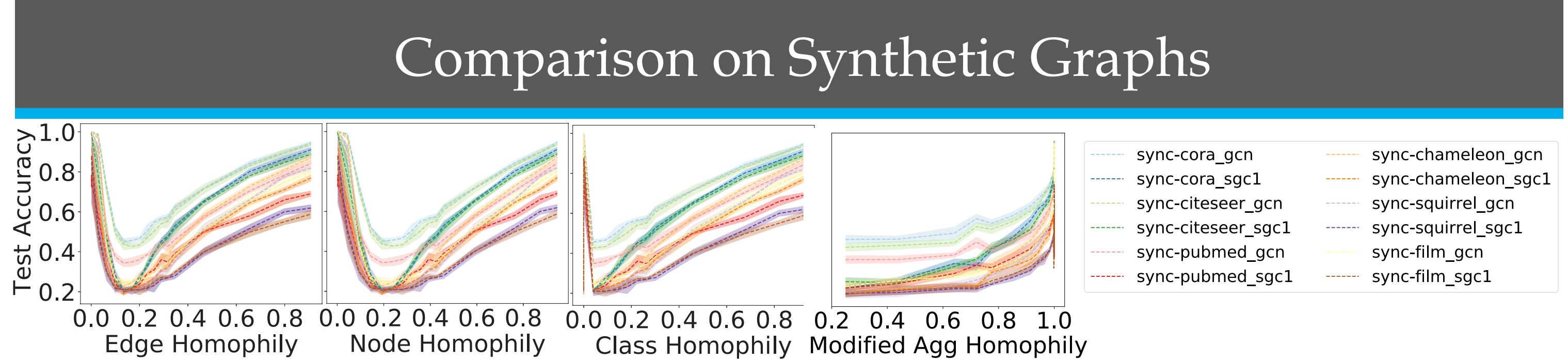
We first define the post-aggregation node similarity matrix as $S(\hat{A}, X) \equiv \hat{A}X(\hat{A}X)^T \in \mathbb{R}^{N \times N}$, where \hat{A} denotes a general aggregation operator and X denotes the feature matrix. $S(\hat{A}, X)$ strongly relates to the gradient of SGC and a large $[S(\hat{A}, X)]_{ij}$ means node i tends to be updated to the same class as node j . This indicates that $S(\hat{A}, X)$ is closely related to the performance of a single layer GNN.

We define the aggregation similarity score $S_{agg}(S(\hat{A}, X))$ as: $\frac{1}{|V|} |\{v | \text{Mean}_u (\{S(\hat{A}, X)_{v,u} | Z_{u,:} = Z_{v,:}\}) \geq \text{Mean}_u (\{S(\hat{A}, X)_{v,u} | Z_{u,:} \neq Z_{v,:}\})\}|$ $S_{agg}(S(\hat{A}, X))$ measures the proportion of nodes v as which the average weights on the set of nodes in the same class (including v) is larger than that in other classes.

Aggregation homophily and its modified version are defined as:

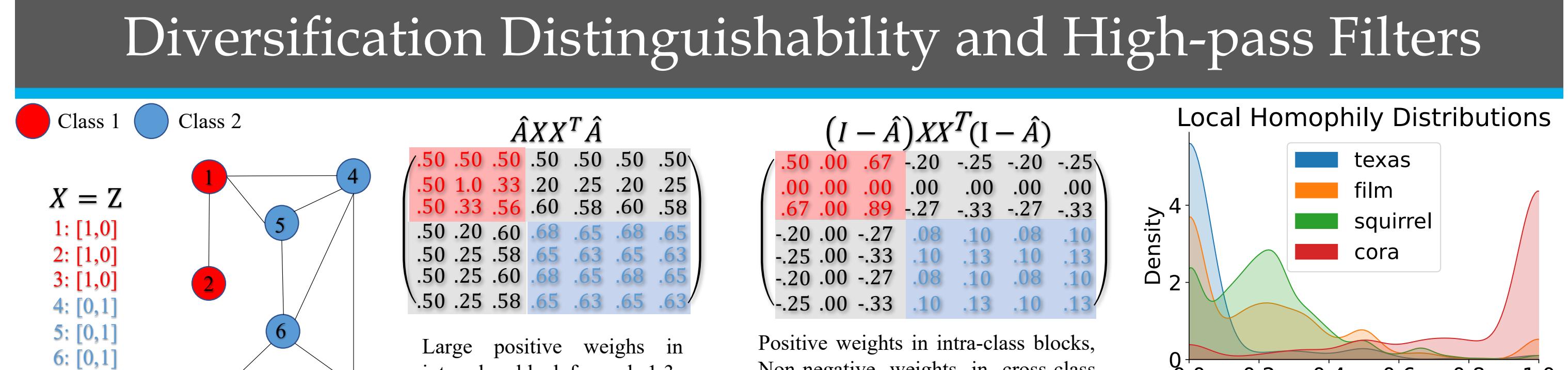
$$H_{agg}(G) = S_{agg}(S(\hat{A}, Z))$$

$$H_{agg}^M(G) = [2 S_{agg}(S(\hat{A}, Z)) - 1]_+$$



In the previous example, when $\hat{A} = \hat{A}_{rw}$, it is easy to see that $H_{agg}(G) = H_{agg}^M(G) = 1$ and other metrics are 0. Thus, this new metric reflects the fact that nodes in classes 1 and 2 are still highly distinguishable after aggregation, while other metrics fail to capture such information and misleadingly give value 0. This shows the advantage of $H_{agg}(G)$ and $H_{agg}^M(G)$, which additionally exploit information from aggregation operator \hat{A} and the similarity matrix.

To comprehensively compare the metrics, we generate synthetic graphs with different homophily levels and evaluate 1-hop SGC and GCN on them. The performance of SGC-1 and GCN is expected to be monotonically increasing if the homophily metric is informative. However, the above figures show that the performance curves under edge, node and class homophily are U-shaped, while for $H_{agg}^M(G)$, it reveals a nearly monotonic curve with a little numerical perturbation around 1. This indicates that $H_{agg}^M(G)$ provides a better indication of the way in which the graph structure affects the performance of SGC-1 and GCN than existing metrics.



Diversification operation, i.e., high-pass(HP) filter, $I - \hat{A}$ is empirically found to be useful for heterophily problem. In the above example, $S(\hat{A}, X)$ assigns large positive weights to pairs of nodes in different classes, which will make nodes hard to be distinguished from each other. However, $S(I - \hat{A}, X)$ assigns positive values to nodes in the same class and negative values to nodes from different classes. This indicates that in some heterophily cases, nodes are still distinguishable through their local surrounding dissimilarities. Thus, we define Diversification Distinguishability (DD) based on $S(I - \hat{A}, X)$.

Definition. Given $S(I - \hat{A}, X)$, a node v is diversification distinguishable if the following two conditions are satisfied at the same time.

1. $\text{Mean}_u (\{S(I - \hat{A}, X)_{v,u} | u \in V, Z_{u,:} = Z_{v,:}\}) \geq 0$; 2. $\text{Mean}_u (\{S(I - \hat{A}, X)_{v,u} | u \in V, Z_{u,:} \neq Z_{v,:}\}) < 0$; Then, graph diversification distinguishability value is defined as

$$DD_{\hat{A}, X}(G) = \frac{1}{|V|} |\{v | v \in V, v \text{ is diversification distinguishable}\}|$$

Theorem: For $C = 2$, suppose $X = Z$, $\hat{A} = \hat{A}_{rw}$, then all nodes are diversification distinguishable and $DD_{\hat{A}, Z}(G) = 1$.

Adaptive Channel Mixing (ACM)

HP filter is proved to be effective for some heterophily cases. Therefore, we propose to include low-pass channel, high-pass channel and identity channel together in each GNN layers instead of using the uni-channel architecture. Besides, different nodes may have different local heterophily situations (see the above figure on homophily distributions) and may have different needs for the information from different channels. So we propose a node-wise channel mixing mechanism to adaptively exploit the channel information. We will use GCN as an example to introduce the ACM framework in matrix form, but the framework can be combined in a similar manner to many different GNNs. The ACM framework includes the following steps:

Step 1. Feature Extraction for Each Channel:

Option 1: $H_L^l = \text{ReLU}(H_{LP}H^{l-1}W_L^{l-1})$, $H_H^l = \text{ReLU}(H_{HP}H^{l-1}W_H^{l-1})$, $H_I^l = \text{ReLU}(H^{l-1}W_I^{l-1})$; Option 2: $H_L^l = H_{LP}\text{ReLU}(H^{l-1}W_L^{l-1})$, $H_H^l = H_{HP}\text{ReLU}(H^{l-1}W_H^{l-1})$, $H_I^l = I\text{ReLU}(H^{l-1}W_I^{l-1})$; $H^0 \in \mathbb{R}^{N \times F_0}, W_L^{l-1}, W_H^{l-1}, W_I^{l-1} \in \mathbb{R}^{F_{l-1} \times F_l}, l = 1, \dots, L$;

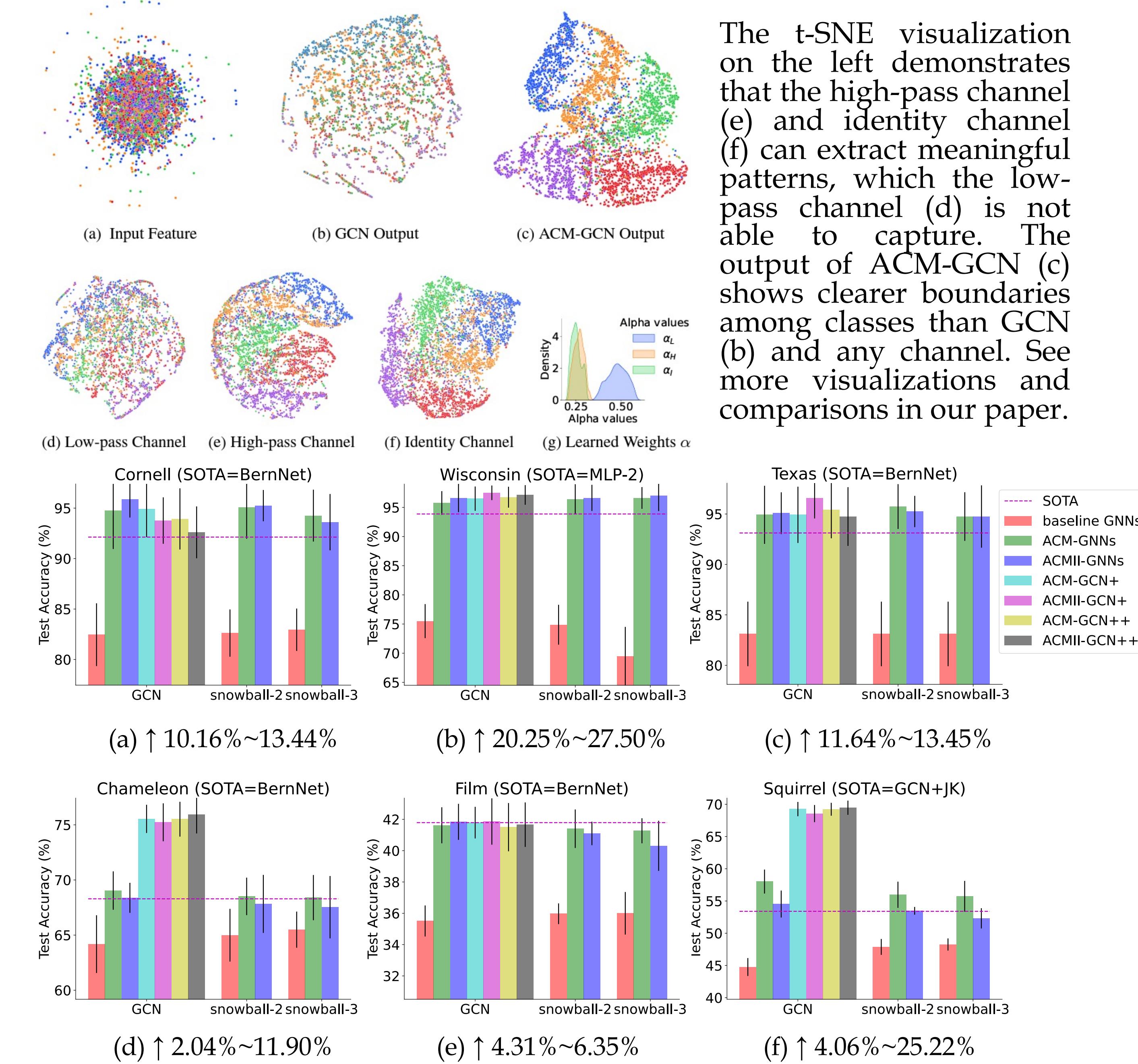
Step 2. Row-wise Feature-based Weight Learning:

$\tilde{\alpha}_L^l = \text{Sigmoid}(H_L^l \tilde{W}_L^l)$, $\tilde{\alpha}_H^l = \text{Sigmoid}(H_H^l \tilde{W}_H^l)$, $\tilde{\alpha}_I^l = \text{Sigmoid}(H_I^l \tilde{W}_I^l)$, $\tilde{W}_L^l, \tilde{W}_H^l, \tilde{W}_I^l \in \mathbb{R}^{F_l \times 1}$; $[\alpha_L^l, \alpha_H^l, \alpha_I^l] = \text{Softmax} \left(\left[\frac{[\tilde{\alpha}_L^l, \tilde{\alpha}_H^l, \tilde{\alpha}_I^l]}{T} \right] W_{\text{Mix}}^l \right) \in \mathbb{R}^{N \times 3}$, $T \in \mathbb{R}$ is temperature, $W_{\text{Mix}}^l \in \mathbb{R}^{3 \times 3}$;

Step 3. Node-wise Adaptive Channel Mixing:

$H^l = \text{ReLU}(\text{diag}(\alpha_L^l)H_L^l + \text{diag}(\alpha_H^l)H_H^l + \text{diag}(\alpha_I^l)H_I^l)$.

Experiments: Node Classification



The t-SNE visualization on the left demonstrates that the high-pass channel (e) and identity channel (f) can extract meaningful patterns, which the low-pass channel (d) is not able to capture. The output of ACM-GCN (c) shows clearer boundaries among classes than GCN (b) and any channel. See more visualizations and comparisons in our paper.

To visualize the performance, we plot the bar charts of the test accuracy of SOTA models, three selected baselines (GCN, snowball-2, snowball-3), their ACM(II) augmented models, ACM(II)-GCN+ and ACM(II)-GCN++ on the 6 most commonly used benchmark heterophily datasets. From the above figures we can see that (1) after being combined with the ACM or ACMII framework, the performance of the three baseline models is **significantly boosted**, by 2.04%~27.50% on all the 6 tasks. The ACM and ACMII in fact achieve SOTA performance. (2) On Cornell, Wisconsin, Texas, Chameleon and Squirrel, the augmented baseline models **significantly outperform the current SOTA models**. Overall, these results suggest that the proposed approach can help GNNs to generalize better on node classification tasks on heterophilic graphs.

See our paper for ablation study, more performance comparisons, details of implementation and discussion of limitations.

Complete the Missing Half: Augmenting Aggregation Filtering with Diversification for Graph Convolutional Networks

Sitao Luan, Mingde Zhao, Chenqing Hua, Xiao-Wen Chang, Doina Precup

Abstract

The core operation of current Graph Neural Networks (GNNs) is the aggregation enabled by the graph Laplacian or message passing, which filters the neighborhood information of nodes. Though effective, in this paper, we show that they are potentially a problematic factor underlying all GNN models for learning on certain datasets, as they force the node representations similar, making the nodes gradually lose their identity and become indistinguishable. Hence, we augment the aggregation operations with their dual, i.e., diversification operators that make the node more distinct and preserve the identity. Such augmentation replaces the aggregation with a two-channel filtering process that, in theory, is beneficial for enriching the node representations. In practice, the proposed two-channel filters can be easily patched on existing GNN methods with diverse training strategies, including spectral and spatial methods. In the experiments, we observe desired characteristics of the models and significant performance boost upon the baselines on 9 node classification tasks.

Measure of Smoothness

Dirichlet Energy can measure the global smoothness of the signal. And the signal energy represents the amount of contents in a graph signal. They are defined as follows.

Dirichlet Energy & Graph Signal Energy The Dirichlet energy of vector block \mathbf{X} and column vector \mathbf{x} defined on \mathcal{G} are separately defined as

$$E_S^G(\mathbf{X}) = \text{tr}(\mathbf{X}^T \mathbf{L} \mathbf{X}), E_S^G(\mathbf{x}) = \mathbf{x}^T \mathbf{L} \mathbf{x}$$

The signal energy of vector block \mathbf{X} and column vector \mathbf{x} defined on \mathcal{G} are separately defined as

$$E_S^G(\mathbf{X}) = \text{tr}(\mathbf{X}^T \mathbf{X}), E_S^G(\mathbf{x}) = \mathbf{x}^T \mathbf{x}$$

$E_S^G(\mathbf{x})$ and $E_S^G(\mathbf{x})$ can be rewritten as

$$E_S^G(\mathbf{x}) = \sum_i \lambda_i \|u_i^T \mathbf{x}\|_2^2, E_S^G(\mathbf{x}) = \sum_i \|u_i^T \mathbf{x}\|_2^2$$

where \mathbf{u}_i and λ_i are eigenvector and eigenvalue of L . Signal \mathbf{x} can be decomposed into smooth and non-smooth components, and the amount of the non-smooth component can be measured by

$$E_{NS}^G(\mathbf{x}) = E_S^G(\mathbf{x}) - E_S^G(\mathbf{x}) = (1 - \lambda_i) \sum_i \|u_i^T \mathbf{x}\|_2^2$$

$E_{NS}^G(\mathbf{x})$ can be negative and a small $E_{NS}^G(\mathbf{x})$ indicates that \mathbf{x} is highly non-smooth. Upon the above analysis, we define $S(\mathbf{x})$ to measure the smoothness of a signal as follows

$$S(\mathbf{x}) = \frac{E_S^G(\mathbf{x})}{E^G(\mathbf{x})}, S(\mathbf{X}) = \frac{E_S^G(\mathbf{X})}{E^G(\mathbf{X})}$$

Signal with a small S -value means it is smooth on \mathcal{G} .

Filterbanks in GNNs: From One Channel to Two

Table 1: Dataset Overview: Network Characteristics and S -values measured by L_{sym}

datasets	Cornell	Wisconsin	Texas	Actor	Chameleon	Squirrel	Cora	CiteSeer	PubMed
Network Info	183	251	183	7600	2277	5201	2708	3327	19717
	295	499	309	33544	36101	217073	5429	4732	44338
	1703	1703	1703	931	2325	2089	1433	3703	500
	5	5	5	5	5	5	7	6	3
S -values	0.904	0.873	0.854	0.901	0.99	0.987	0.862	0.799	0.832
	0.883	0.877	0.909	0.836	0.747	0.782	0.288	0.35	0.501
	-0.021	0.004	0.055	-0.065	-0.243	-0.205	-0.574	-0.449	-0.331

We use blue and red shades to demonstrate the relation between label and feature: the label of the blue shaded datasets is smoother than its feature and red datasets are less smooth.

We measure the smoothness of 9 frequently used benchmark datasets and present the results with the characteristics for each graph in Table 1. It shows that the input features and ground truth labels of different datasets are all mixtures of smooth and non-smooth graph signals but in different proportions. Besides, it illustrates that different tasks have different demands of learning how to smoothen the input signals. For example, in *Cora*, *CiteSeer* and *PubMed*, the ground truth labels are much smoother than the input features, such pattern motivates us to learn how to smoothen the input signals; while in *Wisconsin* and *Texas*, the labels are less smooth than the input features, thus there is no reason that we still learn how to smoothen the input signals. Therefore, to accommodate different situations, we propose that we should learn both smooth and non-smooth components of the input features adaptively instead of merely extracting the smooth part. This motivates us to use filterbanks in GNNs.

Low-pass(LP), High-pass(HP) Graph Filters and Filterbanks The multiplication of L and \mathbf{x} acts as a filtering operation over \mathbf{x} , adjusting the scale of the components of \mathbf{x} in frequency domain. To see this, consider

$$\mathbf{x} = \sum_i \mathbf{u}_i \mathbf{u}_i^T \mathbf{x}, L\mathbf{x} = \sum_i \lambda_i \mathbf{u}_i \mathbf{u}_i^T \mathbf{x},$$

The projection $\mathbf{u}_i \mathbf{u}_i^T \mathbf{x}$ corresponding to a large $|\lambda_i|$ will be amplified, while the one corresponding to a small $|\lambda_i|$ will be suppressed. More specifically, a graph filter that filters out smooth (non-smooth) components is called HP (LP) filter. On the node level, left multiplying HP and LP filters on \mathbf{x} can be understood as diversification and aggregation operations, respectively. For example, if we implement L_{rw} and A_{rw} on the i^{th} node, we have

$$(L_{rw}\mathbf{x})_i = \sum_{j \in N_j} \frac{1}{D_{ij}} (\mathbf{x}_i - \mathbf{x}_j), (A_{rw}\mathbf{x})_i = \sum_{j \in N_j} \frac{1}{D_{ij}} \mathbf{x}_j$$

Intuitively, HP filters depict the differences between one node and its neighbors; While LP filters focus on the similarity within a neighborhood. We believe that these two conjugate components are both indispensable to portray a node. Mathematically, multiplying with LP filter (aggregation) is a linear projection, which will project the features to a fixed subspace. We will lose part of feature information by only using LP filter, and the missing half is the HP component of the learned signals, as $L_{LP} + L_{HP} = I$, which satisfies the perfect reconstruction property.

Filterbank assisted GNNs (FB-GNNs)

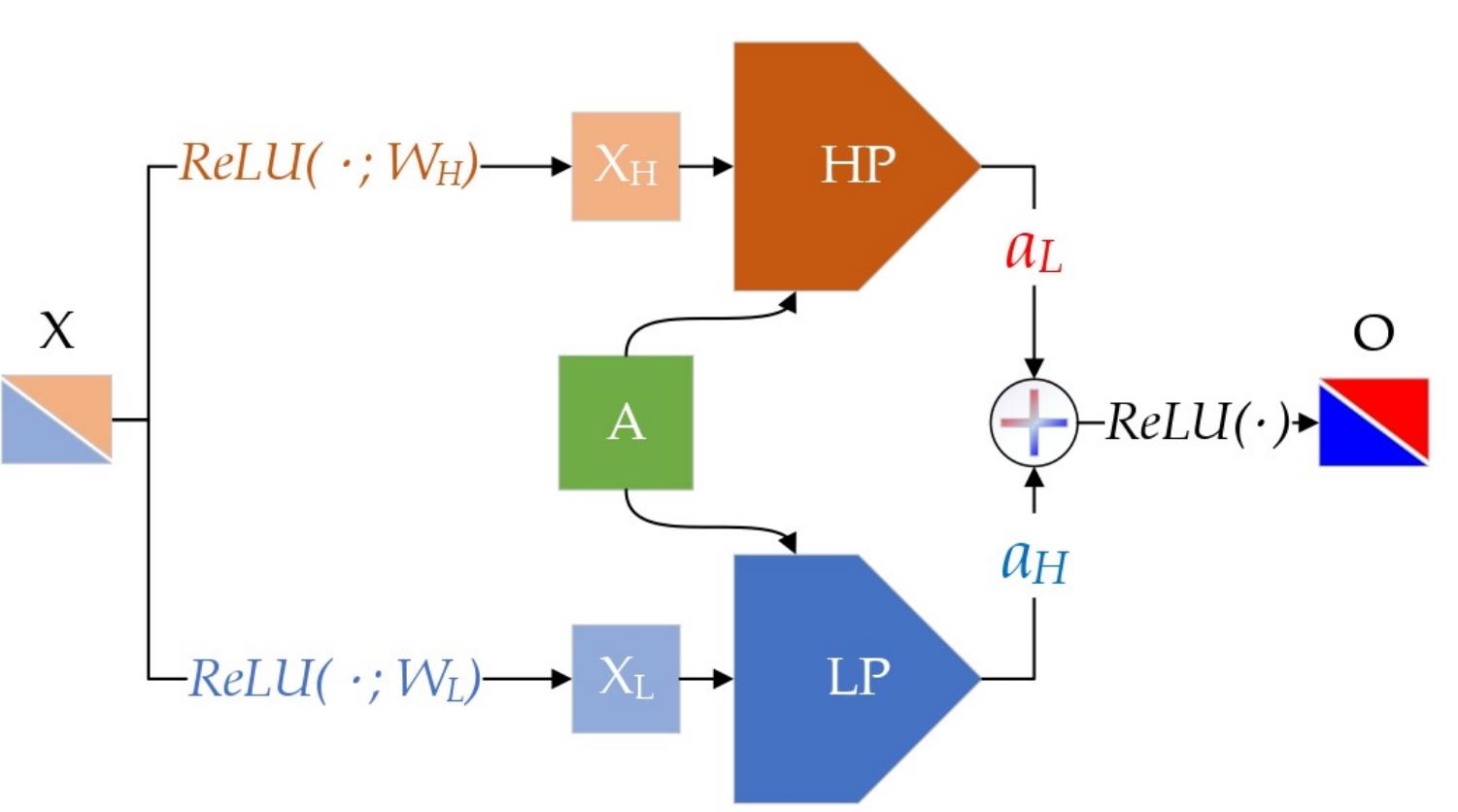
$$H_L^l = L_{LP} \text{ReLU}(H^{l-1} W_L^{l-1})$$

$$H_H^l = L_{HP} \text{ReLU}(H^{l-1} W_H^{l-1})$$

$$H^l = \text{ReLU}(\alpha_L^l \cdot H_L^l + \alpha_H^l \cdot H_H^l)$$

$$l = 1, \dots, n, H^0 = \mathbf{X}, \alpha_L^l, \alpha_H^l \in (0, 1)$$

W_L^{l-1} and W_H^{l-1} are learnable parameter matrices for the non-linear feature extractors focusing on disentangling the smooth and non-smooth information from input H^{l-1} .



Experiments: Node Classification

Shallow GNNs From the results we can see that, our proposed method generally boosts the performance of baseline GNNs in almost all cases, especially when the labels are not much smoother than the input features, considering the S -values in Table 1.

Table 2 shows that the S -values of FB-GCN outputs are closer to the S -values of the ground truth labels (see the absolute differences in the bracket) compared with those of GCN outputs. This indicates that FB-GCN is able to learn better representations which can reconstruct both the smooth and non-smooth part of the ground truth labels.

Model/Datasets	Cornell	Wisconsin	Texas	Actor	Chameleon	Squirrel	Cora	CiteSeer	PubMed
Diff of S -values	-0.023	0.004	0.055	0.065	0.243	0.245	0.578	0.449	-0.330
GAT	70.27	71.35	72.07	73.08	74.06	75.22	76.07	77.32	80.02
FB-GAT	78.36(8.11)	80.06(6.55)	78.66(2.71)	79.08(1.39)	66.05(2.99)	36.06(2.79)	78.75(2.51)	74.79(1.50)	89.00(8.60)
GCN	62.16(9.46)	56.86(10.80)	62.16(10.00)	51.31(13.55)	32.89(6.86)	24.70(7.77)	85.92(6.56)	88.05(6.04)	88.05(5.54)
FB-GCN	64.86(4.05)	72.55(8.41)	70.27(2.70)	31.02(6.01)	67.20(6.30)	49.66(5.52)	85.17(0.24)	76.21(3.09)	88.20(5.01)
GCNII	56.54(8.86)	56.94(2.77)	62.16(2.43)	31.25(9.95)	61.49(1.31)	37.21(1.03)	85.45(0.55)	75.20(1.53)	88.05(1.35)
FB-GCNII	57.30(8.62)	60.68(4.63)	62.11(4.63)	31.45(3.36)	60.76(0.45)	35.21(1.95)	85.45(0.26)	76.33(0.33)	88.05(0.43)
GCNII*	59.64(4.41)	64.87(2.48)	77.78(6.34)	22.34(2.32)	37.36(1.77)	30.61(1.47)	85.61(0.27)	72.83(0.32)	87.90(0.63)
FB-GCNII*	60.11(4.41)	84.67(2.48)	77.78(6.34)	22.34(2.32)	37.36(1.77)	30.61(1.47)	85.61(0.27)	72.83(0.32)	87.90(0.63)

The results are averaged from 10 independent runs. The (brackets) represent the difference of performance brought by patching FB.

Model/Datasets	Cornell	Wisconsin	Texas	Actor	Chameleon	Squirrel	Cora	CiteSeer	PubMed
GCN-B	75.95(4.41)	73.88	71.08	33.77	35.31(6.7)	60.43(3.18)	37.49	75.54	88.62
FB-GCN-B	72.02(7.71)	82.55(8.47)	70.70(7.57)	35.12(1.25)	56.78(1.3)	39.38(1.23)	87.02(5.2)	76.67(6.13)	89.39(1.1)
GCN-C	72.01(11)	76.85(8.47)	80.27(0.81)	32.99(1.38)	54.98(1.27)	36.10(0.83)	88.33(2.02)	76.90(4.41)	89.11(2.1)
FB-GCN-C	76.49(4.41)	76.27(7.41)	76.27(7.34)	29.57(0.81)	54.79(0.25)	36.70(0.8)	87.02(0.57)	76.90(0.08)	89.00(0.08)
GCN-D	76.75(6.79)	82.46(4.12)	78.11(3.74)	34.70(0.81)	41.19(0.5)	36.86(0.56)	76.32(1.26)	70.20(5.3)	89.20(5.26)
FB-GCN-D	76.75(6.79)	82.46(4.12)	78.11(3.74)	34.70(0.81)	41.19(0.5)	36.86(0.56)	76.32(1.26)	70.20(5.3)	89.20(5.26)
GCN-E	74.31(1.19)	80.76(3.72)	84.26(1.35)	34.81(2.2)	59.98(1.12)	42.43(3.54)	76.11(2.02)	76.89(0.34)	89.20(0.26)
FB-GCN-E	75.64(3.41)	81.57(8.04)	80.36(5.13)	34.70(0.95)	57.65(1.2)	39.31(0.29)	87.44(0.57)	77.00(0.1)	89.00(0.32)
GCN-F									

High-Order Pooling For Graph Neural Network With Tensor Decomposition

Chenqing Hua, Guillaume Rabusseau, Jian Tang

ABSTRACT

We propose the Tensorized Graph Neural Network (tGNN), a highly expressive GNN architecture relying on tensor decomposition to model high-order nonlinear node interactions. tGNN leverages the symmetric CP decomposition to efficiently parameterize permutation-invariant multilinear maps for modeling node interactions. Theoretical and empirical analysis on both node and graph classification tasks show the superiority of tGNN over competitive baselines. In particular, tGNN achieves the most solid results on two OGB node classification datasets and one OGB graph classification dataset.

TENSOR

A k-th order tensor $\mathcal{T} \in \mathbb{R}^{N_1 \times N_2 \times \dots \times N_k}$ can simply be seen as a multidimensional array.

The mode-i fibers of \mathcal{T} are the vectors obtained by fixing all indices except the i-th one: $\mathcal{T}_{n_1, n_2, \dots, n_{i-1}, n_{i+1}, \dots, n_k} \in \mathbb{R}^{N_i}$.

The i-th mode matricization of a tensor is the matrix having its mode-i fibers as columns and is denoted by $\mathcal{T}_{(i)}$.

A Rank R CP decomposition factorizes a k-th order tensor $\mathcal{T} \in \mathbb{R}^{N_1 \times N_2 \times \dots \times N_k}$ into the sum of R rank one tensors as

$$\mathcal{T} = \sum_{r=1}^R \mathbf{v}_{1r} \circ \mathbf{v}_{2r} \circ \dots \circ \mathbf{v}_{kr}, \mathbf{v}_{1r} \in \mathbb{R}^{N_1}, \dots, \mathbf{v}_{kr} \in \mathbb{R}^{N_k}.$$

For mode-i product between a tensor $\mathcal{T} \in \mathbb{R}^{N_1 \times N_2 \times \dots \times N_k}$ and a vector $\mathbf{v} \in \mathbb{R}^{N_i}$,

$$(\mathcal{T} \times_i \mathbf{v})_{n_1, \dots, n_{i-1}, n_{i+1}, \dots, n_k} = \sum_{n_i=1}^{N_i} \mathcal{T}_{n_1, \dots, n_k} \mathbf{v}_{n_i}.$$

And the mode-i product with the Kronecker product:

$$\mathcal{T} \times_1 \mathbf{v}_1 \times_2 \mathbf{v}_2 \dots \times_{k-1} \mathbf{v}_{k-1} = \mathcal{T}_{(k)}(\mathbf{v}_{k-1} \otimes \dots \otimes \mathbf{v}_1).$$

CP DECOMPOSITION

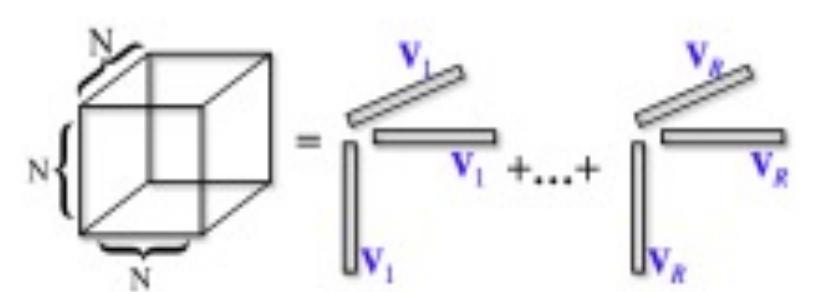


Figure 1: Example of a rank R symmetric CP decomposition of a symmetric 3-order tensor $\mathcal{T} \in \mathbb{R}^{N \times N \times N}$ such that $\mathcal{T} = \sum_{r=1}^R \mathbf{v}_r \circ \mathbf{v}_r \circ \mathbf{v}_r$.

The decomposition vectors, \mathbf{v}_r for $r = 1, \dots, R$, are equal in length, thus can be naturally gathered into factor matrices $\mathbf{M}_1 = [\mathbf{v}_{11}, \dots, \mathbf{v}_{1R}] \in \mathbb{R}^{N_1 \times R}, \mathbf{M}_k = [\mathbf{v}_{k1}, \dots, \mathbf{v}_{kR}] \in \mathbb{R}^{N_k \times R}$.

Using the factor matrices, we denote the CP decomposition of \mathcal{T} as

$$\mathcal{T} = \sum_{r=1}^R \mathbf{v}_{1r} \circ \mathbf{v}_{2r} \circ \dots \circ \mathbf{v}_{kr} = [\mathbf{M}_1, \mathbf{M}_2, \dots, \mathbf{M}_k].$$

The k-th order tensor \mathcal{T} is cubical if all its modes have the same size, i.e., $N_1 = N_2 = \dots = N_k := N$.

A tensor \mathcal{T} is symmetric if it is cubical and is invariant under permutation of its indices:

$$\mathcal{T}_{\mathbf{n}_{\phi(1)}, \dots, \mathbf{n}_{\phi(k)}} = \mathcal{T}_{\mathbf{n}_1, \dots, \mathbf{n}_k} \quad \mathbf{n}_1, \dots, \mathbf{n}_k \in [N]$$

for any permutation $\phi: [k] \rightarrow [k]$.

TENSORIZED GRAPH NEURAL NETWORK

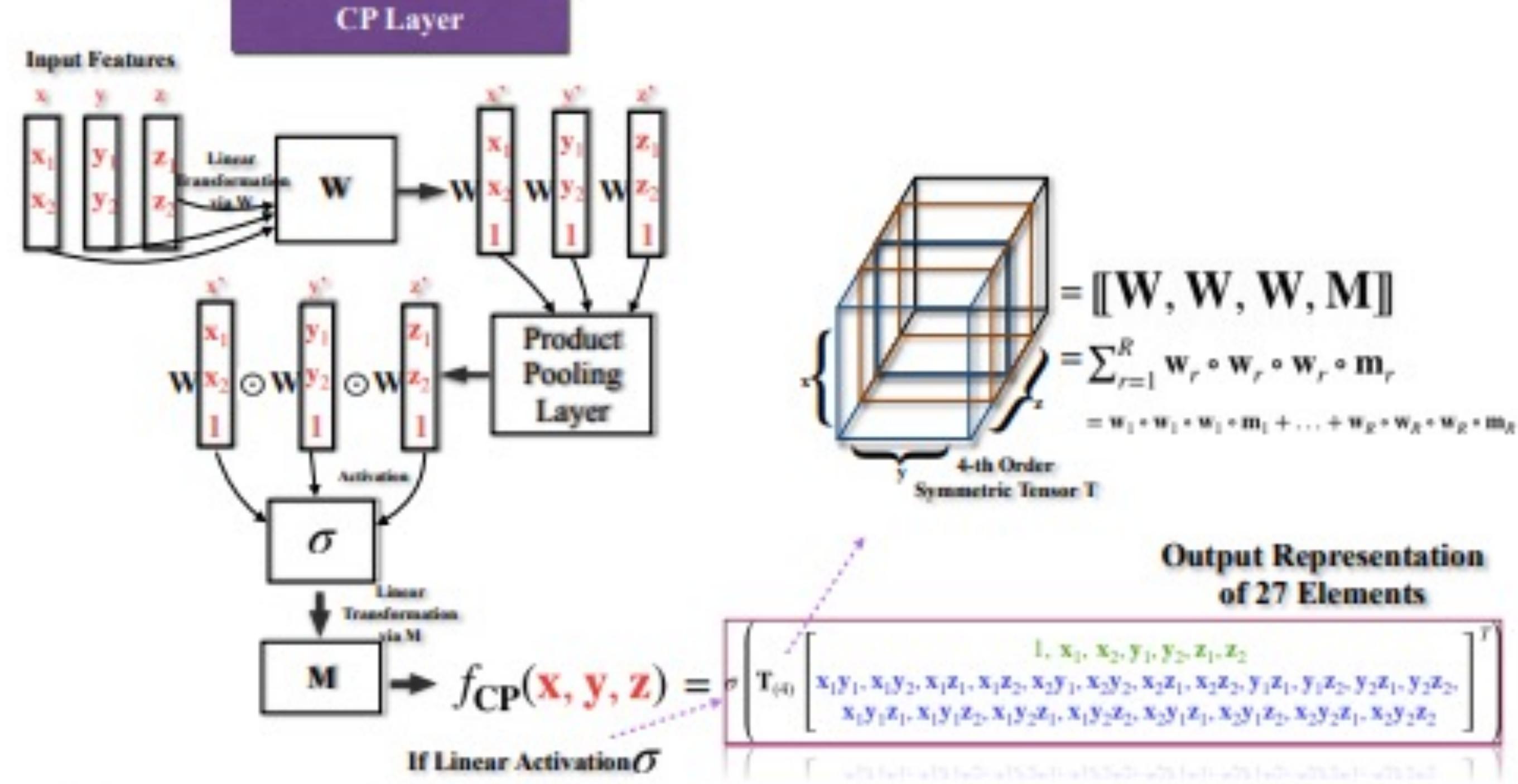


Figure 2: (Left) Sum pooling followed by a FC layer: the output takes individual components of the input into account. (Right) The CP layer can be interpreted as a combination of product pooling with linear layers (with weight matrices \mathbf{W} and \mathbf{M}) and non-linearities. The weight matrices of a CP layer corresponds to a partially symmetric CP decomposition of a weight tensor $\mathcal{T} = [\mathbf{W}, \mathbf{W}, \mathbf{W}, \mathbf{M}]$. It shows that the output of a CP layer takes high-order multiplicative interactions of the inputs' components into account (in contrast with sum pooling that only considers 1st order terms).

Let $\mathcal{T} \in \mathbb{R}^{N \times N \times \dots \times N \times M}$ of order $k+1$ be a tensor which is partially symmetric w.r.t. its first k modes. We can parameterize \mathcal{T} using a rank R partially symmetric CP decomposition: $\mathcal{T} = [\mathbf{W}, \mathbf{W}, \dots, \mathbf{W}, \mathbf{M}], \mathbf{W} \in \mathbb{R}^{N \times R}, \mathbf{M} \in \mathbb{R}^{M \times R}$. Such a tensor naturally defines a map from $(\mathbb{R}^N)^k \rightarrow \mathbb{R}^M$ using contractions over the first k modes:

$$f(\mathbf{x}_1, \dots, \mathbf{x}_k) = \mathcal{T} \times_1 \mathbf{x}_1 \dots \times_{k-1} \mathbf{x}_{k-1} = [\mathbf{W}, \mathbf{W}, \dots, \mathbf{W}, \mathbf{M}] \times_1 \mathbf{x}_1 \dots \times_{k-1} \mathbf{x}_{k-1}.$$

This map satisfies two very important properties for GNNs: it is permutation-invariant (due to the partial symmetry of \mathcal{T}) and its number of parameters is independent of k (due to the partially symmetric CP parameterization).

Definition 1: (CP layer) Given parameter matrices $\mathbf{M} \in \mathbb{R}^{d \times R}, \mathbf{W} \in \mathbb{R}^{(F+1) \times R}$, and activation functions σ, σ' , a rank R CP layer computes the function f_{CP} defined by:

$$f_{CP}(\mathbf{x}_1, \dots, \mathbf{x}_k) = \sigma' \left(\mathbf{M} \left(\sigma \left(\mathbf{W}_1^T \begin{bmatrix} \mathbf{x}_1 \\ 1 \end{bmatrix} \odot \dots \odot \mathbf{W}_k^T \begin{bmatrix} \mathbf{x}_k \\ 1 \end{bmatrix} \right) \right) \right)$$

for any $k \geq 1, \mathbf{x}_1, \dots, \mathbf{x}_k \in \mathbb{R}^F$.

To avoid the numerical instability and counterbalance the bias of high-order CP terms, we use

$$f_{CP}(\mathbf{x}_1, \dots, \mathbf{x}_k) = \sigma' \left(\mathbf{M} \left(\sigma \left(\mathbf{W}_1^T \begin{bmatrix} \mathbf{x}_1 \\ 1 \end{bmatrix} \odot \dots \odot \mathbf{W}_1^T \begin{bmatrix} \mathbf{x}_k \\ 1 \end{bmatrix} \right) \right) \right) + \sigma''(\mathbf{W}_2^T \mathbf{x}_1 + \dots + \mathbf{W}_2^T \mathbf{x}_k),$$

where the first term corresponds to the CP layer and the second one to a standard sum pooling layer (with σ, σ' and σ'' being activation functions).

Theorem 1: The function computed by a CP layer is permutation-invariant. In addition, any permutation-invariant multilinear polynomial $f: (\mathbb{R}^F)^k \rightarrow \mathbb{R}^d$ can be computed by a CP layer (with a linear activation function).

The following theorem shows that CP layers can compute any permutation-invariant multilinear polynomial map. The figure illustrates that the class of functions computed by CP layer subsumes multilinear polynomials (including sum and mean aggregation functions).

Also, the CP layer is strictly more expressive than permutation invariant multilinear polynomials due to the non-linear activation functions.

NODE CLASSIFICATION

Table 1: Results of node-level tasks. **Left Table:** tGNN in comparison with GNN architectures on citation networks. **Right Table:** tGNN in comparison with GNN architectures on OGB datasets.

DATASET MODEL	Cora Acc	Citeseer Acc	Pubmed Acc
GCN	0.8778±0.0096	0.8139±0.0123	0.8890±0.0032
GAT	0.8686±0.0042	0.6720±0.0046	0.8328±0.0012
GraphSAGE	0.8658±0.0026	0.7624±0.0030	0.8658±0.0011
H2GCN	0.8752±0.0061	0.7997±0.0069	0.8778±0.0028
GRPN	0.7951±0.0036	0.6763±0.0038	0.8507±0.0009
APPNP	0.7941±0.0038	0.6859±0.0030	0.8502±0.0009
MixHop	0.6563±0.1131	0.4952±0.1335	0.8704±0.0410
tGNN	0.8808±0.0131	0.8051±0.0192	0.9080±0.0018

DATASET MODEL	PRODUCTS Acc	ARXIV Acc	PROTEINS AUC
MLP	0.6106±0.0008	0.5550±0.0023	0.7204±0.0048
Node2vec	0.7249±0.0010	0.7007±0.0013	0.6881±0.0065
GCN	0.7564±0.0021	0.7174±0.0029	0.7251±0.0035
GraphSAGE	0.7850±0.0016	0.7149±0.0027	0.7768±0.0020
DeeperGCN	0.8098±0.0020	0.7192±0.0016	0.8580±0.0017
Ours	0.8179±0.0054	0.7538±0.0015	0.8255±0.0049

We can observe that tGNN outperforms all classic baselines on Cora, Pubmed, PRODUCTS, and ARXIV, and have slight improvements on the other datasets but underperforms GCN on Citeseer and DeeperGCN on PROTEINS. Regarding the citation networks, tGNN outperforms others on 2 out of 3 datasets. Moreover, on the OGB node datasets, even when tGNN is not ranked first, it is still very competitive (top 3 for all datasets).

GRAPH CLASSIFICATION

Table 2: Results of tGNN on graph-level tasks in comparison with GNN architectures.

DATASET	ZINC	CIFAR10	MNIST	MolHIV
	No edge features MAE	No edge features Acc	No edge features Acc	No edge features AUC
Dwivedi et al.	0.710±0.001	0.560±0.009	0.945±0.003	
GCN	0.469±0.002	0.545±0.001	0.899±0.002	0.761±0.009
GIN	0.408±0.008	0.533±0.037	0.939±0.013	0.756±0.014
DiffPool	0.466±0.006	0.579±0.005	0.950±0.004	
GAT	0.463±0.002	0.655±0.003	0.956±0.001	
MoNet	0.407±0.007	0.534±0.004	0.904±0.005	
GatedGCN	0.422±0.006	0.692±0.003	0.974±0.001	
Corso et al.	0.320±0.032	0.702±0.002	0.972±0.001	0.791±0.013
Le et al.	0.219±0.010	0.727±0.005		0.793±0.012
Beaini et al.	0.219±0.010	0.727±0.005		0.797±0.009
Ours	0.301±0.008	0.684±0.006	0.965±0.002	0.799±0.016

We present tGNN performance on graph property prediction tasks. tGNN achieves state-of-the-art results on MolHIV and has slight improvements on the other three Benchmarking graph datasets.

EFFICIENCY

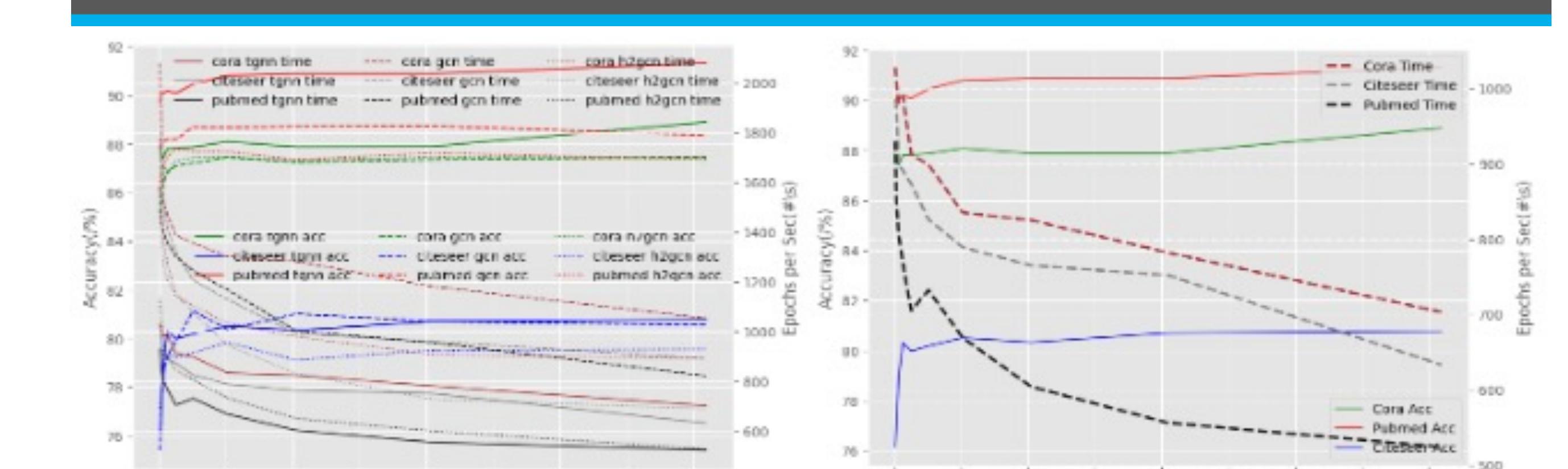


Figure 4: Results of node classification with increasing rank dimension on three citation datasets. **Left Figure:** Left axis shows accuracy, right axis shows #training epochs second, and horizontal axis indicates rank dim of tGNN or hidden dim of baselines. **Right Figure:** Left axis shows accuracy, right axis shows #training epochs second, and horizontal axis indicates rank dim of tGNN.

We can see that the model performance can be improved with higher ranks (i.e., T is more accurately computed as the rank R gets larger), but training time is also increased, thus it is a trade-off between classification accuracy and computation efficiency. In comparison with baselines, tGNN still has marginal improvements with higher ranks while the baselines stop improving with larger hidden dimensions.



Figure 3: Visualization of relations of (permutation-invariant function space) \supseteq {CP function space} \supseteq (permutation-invariant multilinear polynomial space) \supseteq (sum and mean aggregation functions).

When Do Graph Neural Networks Help with Node Classification?

-Investigating the Impact of Homophily Principle on Node Distinguishability

Sitao Luan,¹ Chenqing Hua,¹ Minkai Xu,¹ Qincheng Lu,¹ Jiaqi Zhu,¹ Xiao-Wen Chang,¹ Jie Fu,¹ Jure Leskovec,¹ Doina Precup²

Introduction

In this paper, we point out the deficiency of the current understanding of homophily on Node Distinguishability (ND). To study ND deeply, we propose Contextual Stochastic Block Model for Homophily (CSBM-H) and define two metrics to quantify ND. With the metrics, we visualize and analyze how **graph filters**, **node degree distributions** and **class variances** influence ND, and investigate the combined effect of intra- and inter-class ND. Besides, we verified that, in real-work tasks, the superiority of GNNs is indeed closely related to both intra- and inter-class ND regardless of homophily levels. Grounded in this observation, we propose a new hypothesis-testing based performance metric beyond homophily, which is non-linear, feature-based and can provide statistical threshold value for GNNs' the superiority. Experiments indicate that it is significantly more effective than the existing homophily metrics on revealing the advantage and disadvantage of GNNs on both synthetic and benchmark real-world datasets.

Background & Motivation

The success of GNNs is believed to be rooted in the homophily assumption, which states that connected nodes tend to have similar attributes. On the other hand, the lack of homophily, i.e., heterophily, is considered as the main cause of the inferiority of GNNs on heterophilic graphs. Recently, both empirical and theoretical studies indicate that **the relationship between homophily and GNN performance is more complicated than "homophily wins, heterophily loses"**.

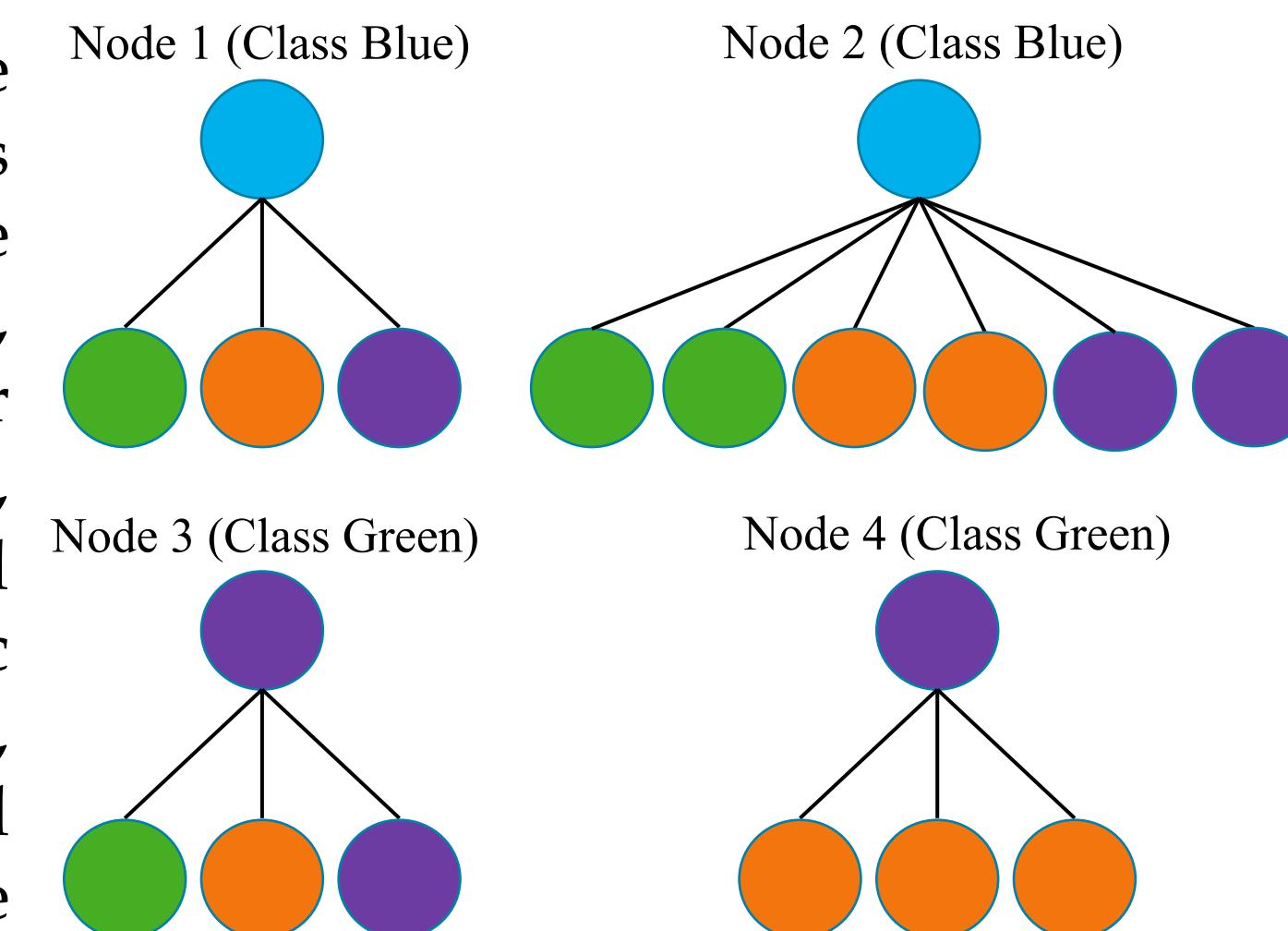
Current understanding: The authors in [1] stated that, as long as nodes within the same class share similar neighborhood patterns, their embeddings will be similar after aggregation. For example, nodes {1,2} are from class blue and both have the same heterophilic neighborhood patterns. As a result, their aggregated features will still be similar, and they can be classified into the same class.

Problem: The above statement is partially true because it only discusses intra-class ND but forgets inter-class ND. For example, node 3 is from class green and has the same neighborhood pattern as {1,2}, which means node 3 and nodes {1,2} will become indistinguishable after aggregation.

Takeaway: This highlights the necessity for careful consideration of both intra- and inter-class ND when evaluating the impact of homophily on the performance of GNNs. An ideal case would be node {1,2,4}, where we have smaller intra-class "distance" than inter-class "distance". We will formulate it.

Quantify Node Distinguishability

CSBM-H The features of 2 classes of nodes, $i \in \mathcal{C}_0$ and $j \in \mathcal{C}_1$, are generated independently, with $\mathbf{x}_i \sim N(\boldsymbol{\mu}_0, \sigma_0^2 \mathbf{I})$ and $\mathbf{x}_j \sim N(\boldsymbol{\mu}_1, \sigma_1^2 \mathbf{I})$, where $\boldsymbol{\mu}_0, \boldsymbol{\mu}_1 \in \mathbb{R}^{F_h}$; d_0, d_1 are node degree for i, j . For $i \in \mathcal{C}_0$, its neighbors are generated by independently sampling from $h \cdot d_0$ intra-class nodes and $(1 - h) \cdot d_0$ inter-class nodes. The neighbors of $j \in \mathcal{C}_1$ are generated in the same way. The centers $(\tilde{\boldsymbol{\mu}}_0, \tilde{\boldsymbol{\mu}}_1, \tilde{\boldsymbol{\mu}}_0^{HP}, \tilde{\boldsymbol{\mu}}_1^{HP})$ and variances $(\tilde{\sigma}_0, \tilde{\sigma}_1, (\tilde{\sigma}_0^{HP})^2, (\tilde{\sigma}_1^{HP})^2)$ of low- and high-pass filtered features can be computed accordingly (See our paper).



Theorem 1 (Optimal Bayes Classifier CL_{Bayes} for CSBM-H) Suppose the prior distribution for \mathbf{x} is $\mathbb{P}(\mathbf{x} \in \mathcal{C}_0) = \mathbb{P}(\mathbf{x} \in \mathcal{C}_1) = 1/2$, then we have

$$CL_{Bayes}(\mathbf{x}) = \begin{cases} 1, & \eta(\mathbf{x}) \geq 0.5 \\ 0, & \eta(\mathbf{x}) < 0.5 \end{cases}, \eta(\mathbf{x}) = \mathbb{P}(z = 1|\mathbf{x}) = \frac{1}{1 + \exp(Q(\mathbf{x}))}$$

where $Q(\mathbf{x}) = \mathbf{a}\mathbf{x}^T + \mathbf{b}^T\mathbf{x} + c$, z_i is the label of node i . See the calculation of a, b, c in the paper.

Bayes Error Rate (BE) for CL_{Bayes} can be written as

$$BE = \mathbb{P}(\mathbf{x} \in \mathcal{C}_0) (1 - \mathbb{P}(CL_{Bayes}(\mathbf{x}) = 0 | \mathbf{x} \in \mathcal{C}_0)) + \mathbb{P}(\mathbf{x} \in \mathcal{C}_1) (1 - \mathbb{P}(CL_{Bayes}(\mathbf{x}) = 1 | \mathbf{x} \in \mathcal{C}_1))$$

Probabilistic Bayes Error (PBE) Since $Q(\mathbf{x})$ in Theorem 1 follows generalized chi-square distribution, BE can be estimated by PBE as follows

$$\text{CDF}_{\tilde{\chi}(\omega_0, F_h, \lambda_0)}(-\xi) + (1 - \text{CDF}_{\tilde{\chi}(\omega_1, F_h, \lambda_1)}(-\xi))$$

$$2$$

See the parameters and the calculation in our paper. PBE can be numerically calculated to show the relation between h and ND precisely. However, it doesn't have an analytic expression, which makes it less explainable.

Generalized Jeffreys Divergence D_{GJ} The KL-divergence based statistical measure offers us a tool to define an explainable ND measure as follows

$$D_{GJ}(P, Q) = \mathbb{P}(\mathbf{x} \in P) \mathbb{E}_{\mathbf{x} \sim P} \left[\ln \frac{P(\mathbf{x})}{Q(\mathbf{x})} \right] + \mathbb{P}(\mathbf{x} \in Q) \mathbb{E}_{\mathbf{x} \sim P} \left[\ln \frac{Q(\mathbf{x})}{P(\mathbf{x})} \right]$$

The negative generalized Jeffreys divergence CSBM-H can be computed by

$$D_{NGJ}(CSBM-H) = -d_X^2 \left(\frac{1}{4\sigma_0^2} + \frac{1}{4\sigma_1^2} \right) - \frac{F_h}{4} [\rho^2 + \frac{1}{\rho^2} - 2]$$

where $d_X^2 = (\mu_0 - \mu_1)^T (\mu_0 - \mu_1)$, $\rho = \sigma_0/\sigma_1$. PBE and D_{GJ} can quantify ND.

Visualization & Analysis

Mid-homophily pitfall: For low-pass (LP) filtered features, the most indistinguishable area is the medium homophily level, instead of the extreme low homophily levels.

3 regimes: low-pass (LP), high-pass (HP) and full-pass (FP) filters dominate different homophily areas. LP: very low and very high homophily intervals (two ends); HP: low to medium homophily levels; FP: medium to high homophily areas.

HP filter: high-pass filter works better in heterophily areas than in homophily areas

Fig 1: $\mu_0 = [-1, 0], \mu_1 = [1, 0], \sigma_0^2 = 1, \sigma_1^2 = 5, d_0 = 5, d_1 = 5$

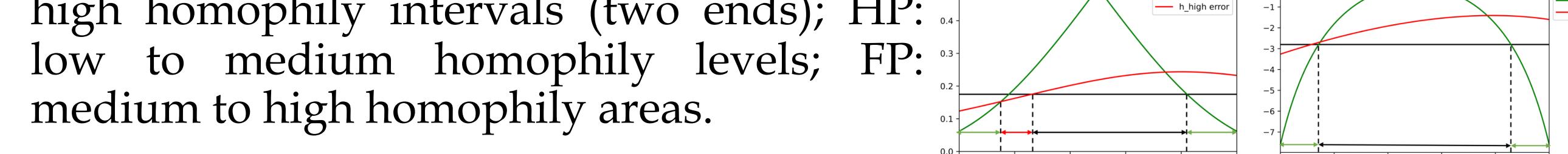


Fig 2: $\mu_0 = [-1, 0], \mu_1 = [1, 0], \sigma_0^2 = 1, \sigma_1^2 = 5, d_0 = 5, d_1 = 5$



Fig 2: $\mu_0 = [-1, 0], \mu_1 = [1, 0], \sigma_0^2 = 1, \sigma_1^2 = 5, d_0 = 5, d_1 = 5$



Fig 2: $\mu_0 = [-1, 0], \mu_1 = [1, 0], \sigma_0^2 = 1, \sigma_1^2 = 5, d_0 = 5, d_1 = 5$



Fig 2: $\mu_0 = [-1, 0], \mu_1 = [1, 0], \sigma_0^2 = 1, \sigma_1^2 = 5, d_0 = 5, d_1 = 5$



Fig 2: $\mu_0 = [-1, 0], \mu_1 = [1, 0], \sigma_0^2 = 1, \sigma_1^2 = 5, d_0 = 5, d_1 = 5$



Fig 2: $\mu_0 = [-1, 0], \mu_1 = [1, 0], \sigma_0^2 = 1, \sigma_1^2 = 5, d_0 = 5, d_1 = 5$



Fig 2: $\mu_0 = [-1, 0], \mu_1 = [1, 0], \sigma_0^2 = 1, \sigma_1^2 = 5, d_0 = 5, d_1 = 5$



Fig 2: $\mu_0 = [-1, 0], \mu_1 = [1, 0], \sigma_0^2 = 1, \sigma_1^2 = 5, d_0 = 5, d_1 = 5$



Fig 2: $\mu_0 = [-1, 0], \mu_1 = [1, 0], \sigma_0^2 = 1, \sigma_1^2 = 5, d_0 = 5, d_1 = 5$



Fig 2: $\mu_0 = [-1, 0], \mu_1 = [1, 0], \sigma_0^2 = 1, \sigma_1^2 = 5, d_0 = 5, d_1 = 5$



Fig 2: $\mu_0 = [-1, 0], \mu_1 = [1, 0], \sigma_0^2 = 1, \sigma_1^2 = 5, d_0 = 5, d_1 = 5$



Fig 2: $\mu_0 = [-1, 0], \mu_1 = [1, 0], \sigma_0^2 = 1, \sigma_1^2 = 5, d_0 = 5, d_1 = 5$



Fig 2: $\mu_0 = [-1, 0], \mu_1 = [1, 0], \sigma_0^2 = 1, \sigma_1^2 = 5, d_0 = 5, d_1 = 5$



Fig 2: $\mu_0 = [-1, 0], \mu_1 = [1, 0], \sigma_0^2 = 1, \sigma_1^2 = 5, d_0 = 5, d_1 = 5$



Fig 2: $\mu_0 = [-1, 0], \mu_1 = [1, 0], \sigma_0^2 = 1, \sigma_1^2 = 5, d_0 = 5, d_1 = 5$



Fig 2: $\mu_0 = [-1, 0], \mu_1 = [1, 0], \sigma_0^2 = 1, \sigma_1^2 = 5, d_0 = 5, d_1 = 5$



Fig 2: $\mu_0 = [-1, 0], \mu_1 = [1, 0], \sigma_0^2 = 1, \sigma_1^2 = 5, d_0 = 5, d_1 = 5$



Fig 2: $\mu_0 = [-1, 0], \mu_1 = [1, 0], \sigma_0^2 = 1, \sigma_1^2 = 5, d_0 = 5, d_1 = 5$



When Do We Need Graph Neural Networks for Node Classification?

Sitao Luan, Chenqing Hua, Qincheng Lu, Jiaqi Zhu, Xiao-Wen Chang, Doina Precup

Abstract

Graph Neural Networks (GNNs) extend basic Neural Networks (NNs) by additionally making use of graph structure based on the relational inductive bias (edge bias), rather than treating the nodes as collections of independent and identically distributed i.i.d. samples. Though GNNs are believed to outperform basic NNs in real-world tasks, it is found that in some cases, GNNs have little performance gain or even underperform graph-agnostic NNs. To identify these cases, based on graph signal processing and statistical hypothesis testing, we propose two measures which analyze the cases in which the edge bias in features and labels does not provide advantages. Based on the measures, a threshold value can be given to predict the potential performance advantages of graph-aware models over graph-agnostic models.

Introduction

Table 1. Accuracy (%) Comparison of Baseline GNNs and MLP

Datasets\Models	MLP Acc	GCN Acc	GAT Acc	GraphSAGE Acc	Baseline Average	Diff(MLP, Edge Homophily)
Cornell	85.14	60.81	59.19	82.97	67.66	17.48 0.3
Wisconsin	87.25	63.73	60.78	87.84	70.78	16.47 0.21
Texas	84.59	61.62	59.73	82.43	67.93	16.66 0.11
Film	36.08	30.98	29.71	35.28	31.99	4.09 0.22
Chameleon	46.21	61.34	61.95	47.32	56.87	-10.66 0.23
Squirrel	29.39	41.86	43.88	30.16	38.63	-9.24 0.22
Cora	74.81	87.32	88.07	85.98	87.12	-12.31 0.81
Citeseer	73.45	76.70	76.42	77.07	76.73	-3.28 0.74
Pubmed	87.86	88.24	87.81	88.59	88.21	-0.35 0.80
DBLP	77.39	85.87	85.89	81.19	84.32	-6.93 0.81
Coauthor CS	93.72	93.91	93.41	94.38	93.90	-0.18 0.81
Coauthor Phy	95.77	96.84	96.32	OOM	96.58	-0.81 0.93
AMZ Comp	83.89	87.03	89.74	83.70	86.82	-2.93 0.78
AMZ Photo	90.87	93.61	94.12	87.97	91.90	-1.03 0.83

Growing evidence shows that GNNs do not always gain advantages over traditional NNs on relational data, e.g., as shown in table 1, MLP outperform baseline GNNs on *Cornell*, *Wisconsin*, *Texas* and *Film* and perform almost the same as baseline GNNs on *PubMed*, *Coauthor CS* and *Coauthor Phy*. This makes us wonder when do we really need GNNs?

Most existing graph filters can be viewed as operators that aggregate node information from its direct neighbors and they are considered as low-pass (LP) filters mainly capturing the low-frequency components of the input, i.e., the locally smooth features across the whole graph. The use of LP graph filters relies on the assumption that nodes tend to share attributes with their neighbors, a tendency called **homophily** that is widely exploited in node classification tasks. GNNs that are built on the homophily assumption learn to assign similar labels to nodes that are closely connected, which corresponds to an assumption of intrinsic smoothness on latent label distribution. We call this kind of relational inductive bias the edge bias. We believe it is a key factor leading to GNNs superiority over NNs.

However, the existing homophily metrics are not appropriate to display the edge bias, e.g. as shown in table 1, MLP does not necessarily outperform baseline GNNs on some low homophily datasets (*Chameleon* and *Squirrel*)

and does not significantly underperform baseline GNNs on some high homophily datasets (*PubMed*, *Coauthor CS*, *Coauthor Phy* and *AMZ Photo*). Thus, a metric that is able to indicate whether or not the graph-aware models can outperform graph-agnostic models is needed.

Measuring the Effect of Edge Bias

Normalized Total Variation (NTV): it measures the distance between node features and its aggregated neighborhood features

$$NTV(X) = \frac{\|X - \hat{A}X\|_2^2}{\|X\|_2^2}$$

Normalized Smoothness Value (NSV): it measures how different is the average pairwise attribute distance of connected nodes from that of unconnected nodes. First, we have the total pairwise attribute distance of connected nodes (Dirichlet energy)

$$E_D^G(X) = \text{Trace}(X^\top LX) = \sum_{i \leftrightarrow j} \|X_{i,:} - X_{j,:}\|_2^2$$

The total pairwise distance of unconnected nodes can be derived from the Laplacian $L^C = D^C - A^C = NI - \mathbf{1}^\top \mathbf{1} - L$ of the complementary graph \mathcal{G}^C

$$E_D^{G^C}(X) = \text{Trace}(X^\top L^C X) = \sum_{i \leftrightarrow j} \|X_{i,:} - X_{j,:}\|_2^2$$

Then, the average pairwise distance of connected and unconnected nodes can be calculated by normalizing $E_D^G(X)$ and $E_D^{G^C}(X)$

$$E_N^G(X) = \frac{E_D^G(X)}{2|\mathcal{E}|}, E_N^{G^C}(X) = \frac{E_D^{G^C}(X)}{N(N-1) - 2|\mathcal{E}|}$$

Then, the Normalized Smoothness Value (NSV) is defined as

$$NSV(X) = \frac{E_N^G(X)}{E_N^G(X) + E_N^{G^C}(X)}$$

Analysis $0 \leq NSV(X) \leq 1$ and it can help to interpret the edge bias: (1) For labels Z , $NSV(Z) \geq 0.5$ means that the proportion of connected nodes that share different labels is larger than that of unconnected nodes, which implies that edge bias is harmful for the prediction of Z and the homophily assumption is invalid; (2) For features X , $NSV(X) \geq 0.5$ means that the average pairwise feature distance of connected nodes is greater than that of unconnected nodes, which suggests that the feature is non-smooth. On the contrary, small $NSV(Z)$ and $NSV(X)$ indicates that the homophily assumption holds and the edge bias is potentially beneficial.

The above analysis raises another question: how much does NSV deviating from 0.5 or what is the exact NSV to indicate the edge bias is statistically significantly beneficial or harmful?

Hypothesis Testing for Edge Bias

Consider the following distributions of labels and features. For labels Z ,

$$P_1 = \mathbb{P}(Z_{i,:} \neq Z_{j,:} | e_{ij} \in \mathcal{E}) = \text{The proportion of connected nodes that share different labels};$$

$$P_2 = \mathbb{P}(Z_{i,:} \neq Z_{j,:} | e_{ij} \notin \mathcal{E}) = \text{The proportion of unconnected nodes that share different labels}.$$

For features X ,

$$D_1 = \|X_{i,:} - X_{j,:}\|_2^2 | e_{ij} \in \mathcal{E} = \text{Distribution of pairwise feature distance of connected nodes};$$

$$D_2 = \|X_{i,:} - X_{j,:}\|_2^2 | e_{ij} \notin \mathcal{E} = \text{Distribution of pairwise feature distance of unconnected nodes}.$$

Suppose P_1, P_2, D_1, D_2 follow

$$P_1 \sim \text{Binom}(n_1, p_1), P_2 \sim \text{Binom}(n_2, p_2); D_1 \sim N(d_1, \sigma_1^2), D_2 \sim N(d_2, \sigma_2^2)$$

Consider the hypotheses for labels

$$H_0^L: p_1 = p_2; H_1^L: p_1 \neq p_2; H_2^L: p_1 \geq p_2; H_3^L: p_1 \leq p_2$$

and hypotheses for features

$$H_0^F: d_1 = d_2; H_1^F: d_1 \neq d_2; H_2^F: d_1 \geq d_2; H_3^F: d_1 \leq d_2$$

To conduct the hypothesis tests, we use Welch's t-test for features and χ^2 test for labels. We can see $E_N^G(Z)$ and $E_N^{G^C}(Z)$ are sample estimation of the p_1, p_2 for label Z ; $E_N^G(X)$ and $E_N^{G^C}(X)$ are sample estimation of d_1, d_2 for X . Thus, the p-values of hypothesis tests can suggest if NSV statistically significantly deviates from 0.5.

Results and Analysis

Analysis For feature X :

- p-value (H_0^F vs. H_1^F): > 0.05 , H_0^F holds, feature is non-smooth; ≤ 0.05 , TBD.
- p-value (H_0^F vs. H_2^F): ≤ 0.05 , feature is significantly non-smooth.
- p-value (H_0^F vs. H_3^F): ≤ 0.05 , feature is statistically significantly smooth.

For label Z :

- p-value (H_0^L vs. H_1^L): > 0.05 , H_0^L holds, label is non-smooth; ≤ 0.05 , TBD.
- p-value (H_0^L vs. H_2^L): ≤ 0.05 , label is statistically significantly non-smooth.
- p-value (H_0^L vs. H_3^L): ≤ 0.05 , label is statistically significantly smooth.

Table 2. Statistics of Datasets and the Performance Differences

Datasets\Measures	Features						Labels						Baseline Average
	NTV	NSV	H_0^F vs H_1^F	H_0^F vs H_2^F	H_0^F vs H_3^F	NTV	NSV	H_0^L vs H_1^L	H_0^L vs H_2^L	H_0^L vs H_3^L	MLP		
Cornell	0.33	0.48	0.00	1.00	0.00	0.33	0.53	0.0003	0.00	1.00	17.48	-	
Texas	0.33	0.48	0.00	1.00	0.00	0.42	0.60	0.00	0.00	1.00	16.66	-	
Wisconsin	0.38	0.51	0.72	0.36	0.64	0.40	0.55	0.00	0.00	1.00	16.47	-	
Film	0.39	0.50	0.19	0.90	0.10	0.37	0.50	0.05	0.97	0.03	4.09	-	
Coauthor CS	0.36	0.36	0.00	1.00	0.00	0.19	0.18	0.00	1.00	0.00	0.18	-	
Pubmed	0.33	0.44	0.00	1.00	0.00	0.25	0.24	0.00	1.00	0.00	0.35	-	
Coauthor Phy	0.35	0.36	0.00	1.00	0.00	0.16	0.09	0.00	1.00	0.00	0.81	-	
AMZ Photo	0.41	0.39	0.00	1.00	0.00	0.23	0.17	0.00	1.00	0.00	1.03	-	
AMZ Comp	0.41	0.38	0.00	1.00	0.00	0.25	0.22	0.00	1.00	0.00	2.93	-	
Citeseer	0.35	0.45	0.00	1.00	0.00	0.22	0.24	0.00	1.00	0.00	3.28	-	
DBLP	0.37	0.46	0.00	1.00	0.00	0.21	0.20	0.00	1.00	0.00	6.93	-	
Squirrel	0.47	0.54	0.00	1.00	0.00	0.44	0.49	0.00	1.00	0.00	9.24	-	
Chameleon	0.45	0.45	0.00	1.00	0.00	0.45							