

# *How Powerful are Graph Neural Networks?*

Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka

ICLR 2019 (oral)

Yizhu Jiao

2019/06/04

# Introduction

- Tasks on graph
  - Graph classification
  - Node classification
  - Link prediction
- Currently used framework in GNN
  - neighborhood aggregation (message passing)

# Graph Neural Networks (GNN)

- Iteratively update the representation of nodes

$$a_v^{(k)} = \text{AGGREGATE}^{(k)} \left( \left\{ h_u^{(k-1)} : u \in \mathcal{N}(v) \right\} \right)$$

$$h_v^{(k)} = \text{COMBINE}^{(k)} \left( h_v^{(k-1)}, a_v^{(k)} \right)$$

- Get the entire graph's representation

$$h_G = \text{READOUT} \left( \left\{ h_v^{(K)} \mid v \in G \right\} \right)$$

# GNN variants

- GraphSAGE (Hamilton et al., 2017)
  - AGGREGATE

$$a_v^{(k)} = \text{MAX} \left( \left\{ \text{ReLU} \left( W \cdot h_u^{(k-1)} \right), \forall u \in \mathcal{N}(v) \right\} \right)$$

- COMBINE

$$W \cdot \left[ h_v^{(k-1)} \parallel a_v^{(k)} \right]$$

# GNN variants

- Graph Convolutional Networks (GCN)
  - AGGREGATE and COMBINE are integrated

$$h_v^{(k)} = \text{ReLU} \left( W \cdot \text{MEAN} \left\{ h_u^{(k-1)}, \forall u \in \mathcal{N}(v) \cup \{v\} \right\} \right)$$

# Motivations

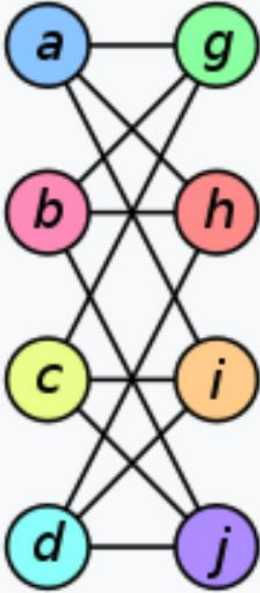
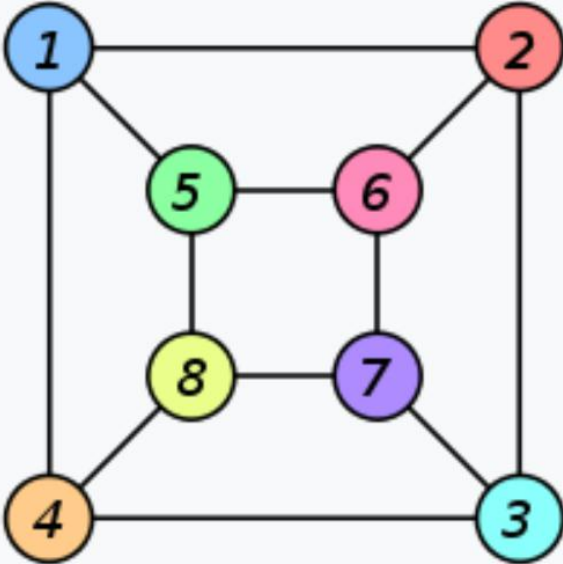
- Design of new GNNs
  - mostly based on empirical intuition and experimental trial-and-error
- How to analyze the representational power of GNNs ?
- How to design a more powerful GNN ?

# Theoretical Framework

- Base on neighborhood aggregation framework
- Base on graph classification
  - distinguish different graph structures
    - imply solving graph isomorphism
  - capture different graph's structural similarity

# Theoretical Framework

- Graph isomorphism

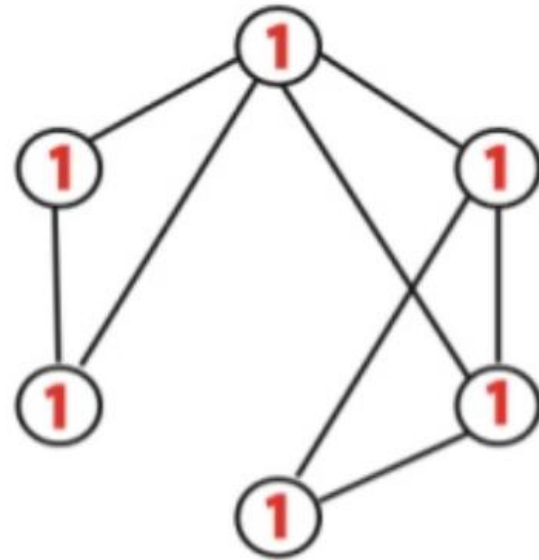
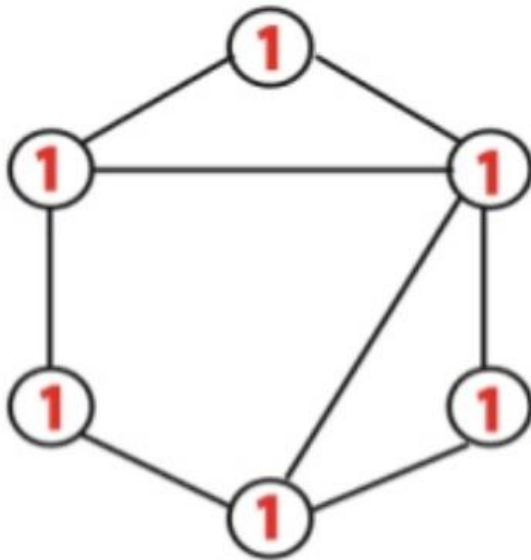
图 $G$	图 $H$	从图 $G$ 到图 $H$ 的同构映射 $\sigma$
		$\begin{aligned}\sigma(a) &= 1 \\ \sigma(b) &= 6 \\ \sigma(c) &= 8 \\ \sigma(d) &= 3 \\ \sigma(g) &= 5 \\ \sigma(h) &= 2 \\ \sigma(i) &= 4 \\ \sigma(j) &= 7\end{aligned}$



# Theoretical Framework

- weaker criterion: Weisfeiler-Lehman (WL) test

Given two graphs  $G$  and  $G'$



# Theoretical Framework

- WL graph isomorphism test

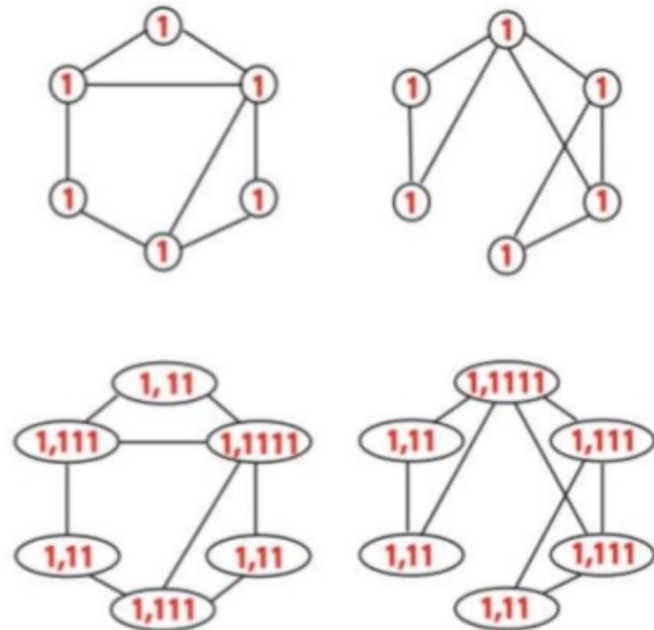
## WL Algorithm :iteration 1

---

Each Iteration of WL test comprises of following steps:-

1. Multiset label determination and sorting

- $O(m)$  via Bucket Sort



# Theoretical Framework

- WL graph isomorphism test

## WL Algorithm :iteration 1

---

Each Iteration of WL test comprises of following steps:-

1.Multiset label determination and sorting

- $O(m)$  via Bucket Sort

2.label compression

- $O(m)$  via Radix Sort

1,11	1,11	1,111
1,11	1,11	1,111
1,11	1,111	1,1111
1,11	1,111	1,1111

1,11	→	2
1,111	→	3
1,1111	→	4

# Theoretical Framework

- WL graph isomorphism test

## WL Algorithm :iteration 1

Each Iteration of WL test comprises of following steps:-

1. Multiset label determination and sorting

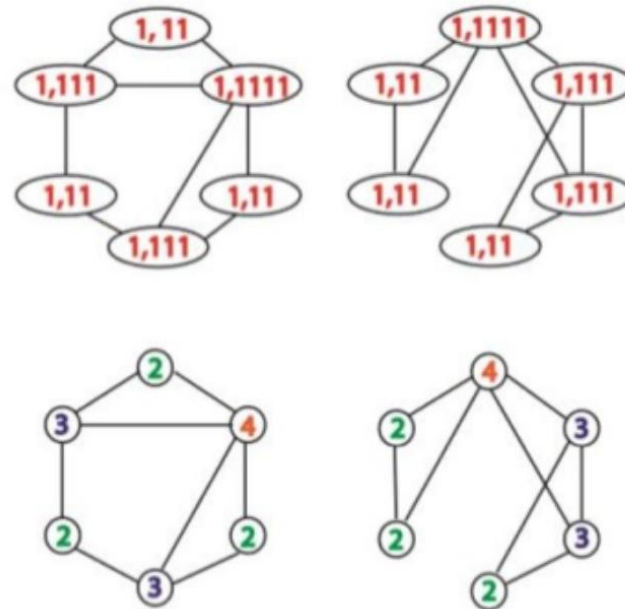
- $O(m)$  via Bucket Sort

2. label compression

- $O(m)$  via Radix Sort

3. Relabeling

- $O(n)$



# Theoretical Framework

- WL graph isomorphism test

## WL Algorithm :iteration 1

---

Each Iteration of WL test comprises of following steps:-

1.Multiset label determination and sorting

- $O(m)$  via Bucket Sort

2.label compression

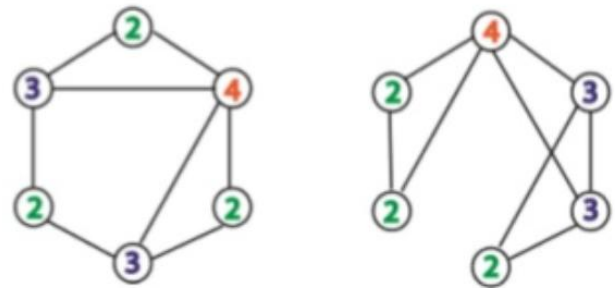
- $O(m)$  via Radix Sort

3. Relabeling

- $O(n)$

4.Are the labels of  $G$  and  $G'$  identical?

Yes, continue.



# Theoretical Framework

- WL graph isomorphism test

## WL Algorithm :iteration 2

---

Each Iteration of WL test comprises of following steps:-

1.Multiset label determination and sorting

◦  $O(m)$  via Bucket Sort

2.label compression

◦  $O(m)$  via Radix Sort

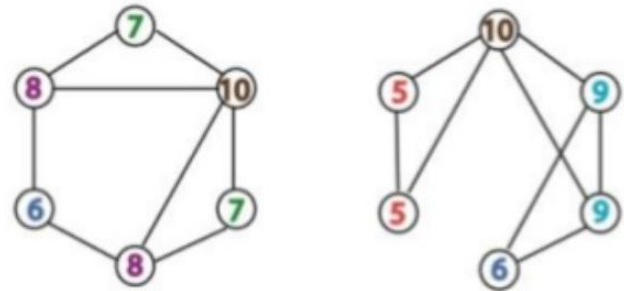
3. Relabeling

◦  $O(n)$

4.Are the labels of  $G$  and  $G'$  identical?

NO, output YES.

5.complexity  $O(hm)$  for  $h$  iteration



# Theoretical Framework

- Multiset
  - the feature vectors of a node's neighbors

$$X = (S, m)$$

# Theoretical Framework

**Lemma 2.** *Let  $G_1$  and  $G_2$  be any non-isomorphic graphs. If a graph neural network  $\mathcal{A} : \mathcal{G} \rightarrow \mathbb{R}^d$  following the neighborhood aggregation scheme maps  $G_1$  and  $G_2$  to different embeddings, the Weisfeiler-Lehman graph isomorphism test also decides  $G_1$  and  $G_2$  are not isomorphic.*

- GNNs are at most as powerful as the WL test in distinguishing graph structures



# Theoretical Framework

**Theorem 3.** *Let  $\mathcal{A} : \mathcal{G} \rightarrow \mathbb{R}^d$  be a GNN following the neighborhood aggregation scheme. With sufficient iterations,  $\mathcal{A}$  maps any graphs  $G_1$  and  $G_2$  that the Weisfeiler-Lehman test of isomorphism decides as non-isomorphic, to different embeddings if the following conditions hold:*

*a)  $\mathcal{A}$  aggregates and updates node features iteratively with*

$$h_v^{(k)} = \phi \left( h_v^{(k-1)}, f \left( \left\{ h_u^{(k-1)} : u \in \mathcal{N}(v) \right\} \right) \right),$$

*where the functions  $f$ , which operates on multisets, and  $\phi$  are injective.*

*b)  $\mathcal{A}$ 's graph-level readout, which operates on the multiset of node features  $\{h_v^{(k)}\}$ , is injective.*

- GNN maps different graph structures to different embedding if AGGREGATE, COMBINE, READOUT are injective

# Graph Isomorphism Network (GIN)

**Corollary 6.** Assume  $\mathcal{X}$  is countable. There exists a function  $f : \mathcal{X} \rightarrow \mathbb{R}^n$  so that for infinitely many choices of  $\epsilon$ , including all irrational numbers,  $h(c, X) = (1 + \epsilon) \cdot f(c) + \sum_{x \in X} f(x)$  is unique for each pair  $(c, X)$ , where  $c \in \mathcal{X}$  and  $X \subset \mathcal{X}$  is a finite multiset. Moreover, any function  $g$  over such pairs can be decomposed as  $g(c, X) = \phi((1 + \epsilon) \cdot f(c) + \sum_{x \in X} f(x))$  for some function  $\phi$ .

- $g(c, X) = \phi((1 + \epsilon) \cdot f(c) + \sum_{x \in X} f(x))$
- model  $f^{(k+1)} \circ \phi^{(k)}$  with one MLP

# Graph Isomorphism Network (GIN)

- AGGREGATE and COMBINE

$$h_v^{(k)} = \text{MLP}^{(k)} \left( \left(1 + \epsilon^{(k)}\right) \cdot h_v^{(k-1)} + \sum_{u \in \mathcal{N}(v)} h_u^{(k-1)} \right)$$

- READOUT

$$h_G = \text{CONCAT} \left( \text{READOUT} \left( \left\{ h_v^{(k)} \mid v \in G \right\} \right) \mid k = 0, 1, \dots, K \right)$$

- GIN provably generalizes the WL test

# Analysis on less powerful GNNs

- Ablation on the aggregator in

$$h_v^{(k)} = \text{MLP}^{(k)} \left( \left(1 + \epsilon^{(k)}\right) \cdot h_v^{(k-1)} + \sum_{u \in \mathcal{N}(v)} h_u^{(k-1)} \right)$$

- Two aspects
  - 1-layer perceptrons instead of MLPs
  - Mean or Max-pooling instead of the sum
    - Mean -> GCN
    - Max -> GraphSAGE

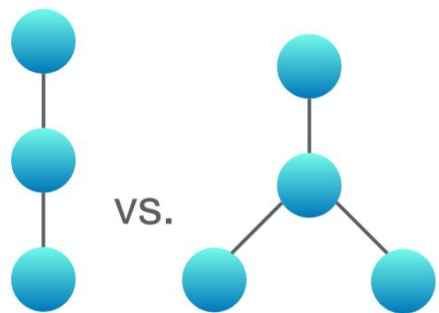
# Analysis on less powerful GNNs

**Lemma 7.** *There exist finite multisets  $X_1 \neq X_2$  so that for any linear mapping  $W$ ,  $\sum_{x \in X_1} \text{ReLU}(Wx) = \sum_{x \in X_2} \text{ReLU}(Wx)$ .*

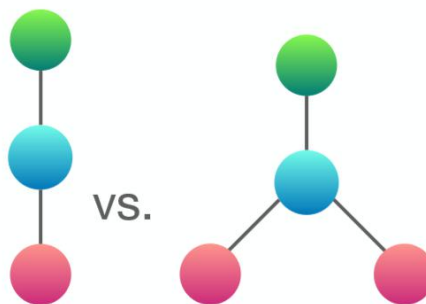
- Linear mapping + ReLU is not sufficient
- Linear mapping + bias + ReLU can distinguish to some degree
  - May not adequately capture structural similarity
  - Difficult for simple classifiers to fit

# Analysis on less powerful GNNs

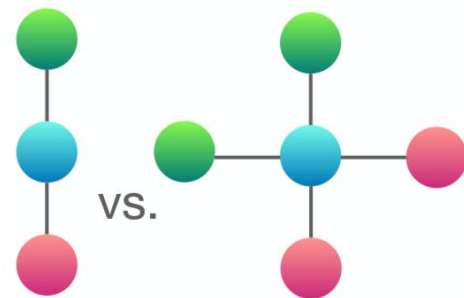
- Structures that confuse mean and max-pooling



(a) Mean and Max both fail



(b) Max fails



(c) Mean and Max both fail

- $f(a)$  and  $2 \cdot f(a)$  are different

# Analysis on less powerful GNNs

- Mean-pooling learns distributions

$$X_1 = (S, m) \quad X_2 = (S, k \cdot m)$$

- Mean aggregator is as powerful as the sum aggregator if node features are diverse and rarely repeat

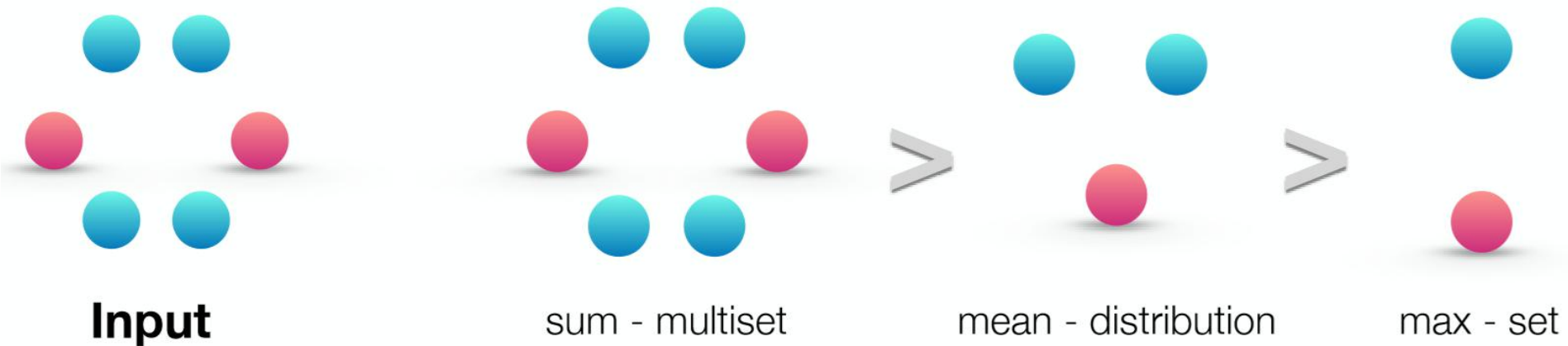
# Analysis on less powerful GNNs

- Max-pooling learns sets with distinct elements
  - Treat a multiset as a set
  - Capture graph skeleton
- Max aggregator may be suitable for tasks where it is important to identify representative elements



# Analysis on less powerful GNNs

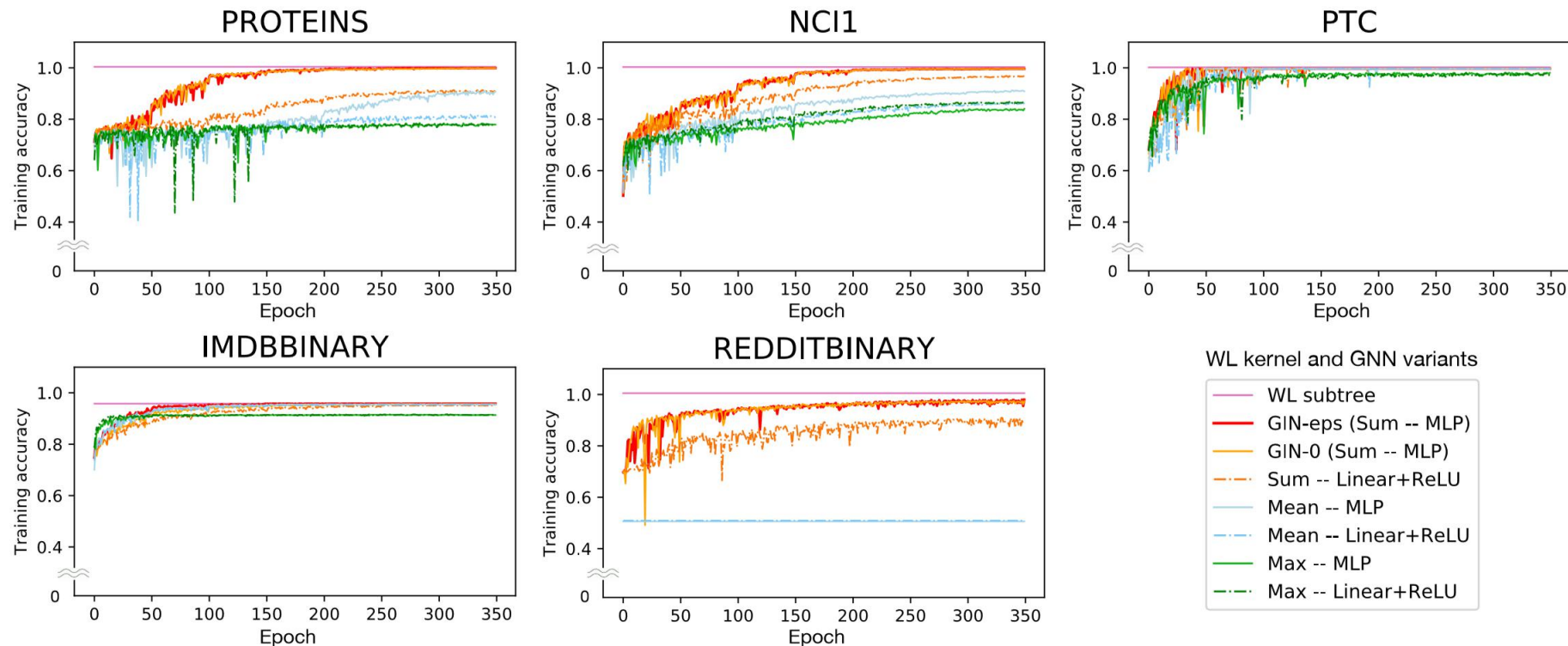
- Ranking by expressive power



# Analysis on less powerful GNNs

- Other neighbor aggregations (not message-passing-based)
  - Graph Attention Networks (GAT)
    - Self-attention
  - LSTM pooling (Hamilton et al., 2017a; Murphy et al., 2018)

# Experiments



Training set performance

# Experiments

Datasets	Datasets	IMDB-B	IMDB-M	RDT-B	RDT-M5K	COLLAB	MUTAG	PROTEINS	PTC	NCI1
	# graphs	1000	1500	2000	5000	5000	188	1113	344	4110
	# classes	2	3	2	5	3	2	2	2	2
	Avg # nodes	19.8	13.0	429.6	508.5	74.5	17.9	39.1	25.5	29.8
Baselines	WL subtree	$73.8 \pm 3.9$	$50.9 \pm 3.8$	$81.0 \pm 3.1$	$52.5 \pm 2.1$	$78.9 \pm 1.9$	$90.4 \pm 5.7$	$75.0 \pm 3.1$	$59.9 \pm 4.3$	<b><math>86.0 \pm 1.8^*</math></b>
	DCNN	49.1	33.5	–	–	52.1	67.0	61.3	56.6	62.6
	PATCHYSAN	$71.0 \pm 2.2$	$45.2 \pm 2.8$	$86.3 \pm 1.6$	$49.1 \pm 0.7$	$72.6 \pm 2.2$	<b><math>92.6 \pm 4.2^*</math></b>	$75.9 \pm 2.8$	$60.0 \pm 4.8$	$78.6 \pm 1.9$
	DGCNN	70.0	47.8	–	–	73.7	85.8	75.5	58.6	74.4
	AWL	$74.5 \pm 5.9$	$51.5 \pm 3.6$	$87.9 \pm 2.5$	$54.7 \pm 2.9$	$73.9 \pm 1.9$	$87.9 \pm 9.8$	–	–	–
GNN variants	GIN- $\epsilon$ (SUM-MLP)	<b><math>74.3 \pm 5.1</math></b>	<b><math>52.1 \pm 3.6</math></b>	<b><math>92.2 \pm 2.3</math></b>	<b><math>57.0 \pm 1.7</math></b>	<b><math>80.1 \pm 1.9</math></b>	<b><math>89.0 \pm 6.0</math></b>	<b><math>75.9 \pm 3.8</math></b>	$63.7 \pm 8.2$	<b><math>82.7 \pm 1.6</math></b>
	GIN-0 (SUM-MLP)	<b><math>75.1 \pm 5.1</math></b>	<b><math>52.3 \pm 2.8</math></b>	<b><math>92.4 \pm 2.5</math></b>	<b><math>57.5 \pm 1.5</math></b>	<b><math>80.2 \pm 1.9</math></b>	<b><math>89.4 \pm 5.6</math></b>	<b><math>76.2 \pm 2.8</math></b>	<b><math>64.6 \pm 7.0</math></b>	<b><math>82.7 \pm 1.7</math></b>
	SUM-1-LAYER	$74.1 \pm 5.0$	<b><math>52.2 \pm 2.4</math></b>	$90.0 \pm 2.7$	$55.1 \pm 1.6$	<b><math>80.6 \pm 1.9</math></b>	<b><math>90.0 \pm 8.8</math></b>	<b><math>76.2 \pm 2.6</math></b>	$63.1 \pm 5.7$	$82.0 \pm 1.5$
	MEAN-MLP	$73.7 \pm 3.7$	<b><math>52.3 \pm 3.1</math></b>	$50.0 \pm 0.0^\dagger$ ( $71.2 \pm 4.6$ )	$20.0 \pm 0.0^\dagger$ ( $41.3 \pm 2.1$ )	$79.2 \pm 2.3$	$83.5 \pm 6.3$	$75.5 \pm 3.4$	<b><math>66.6 \pm 6.9</math></b>	$80.9 \pm 1.8$
	MEAN-1-LAYER	$74.0 \pm 3.4$	$51.9 \pm 3.8$	$50.0 \pm 0.0^\dagger$ ( $69.7 \pm 3.2$ )	$20.0 \pm 0.0^\dagger$ ( $39.7 \pm 2.4$ )	$79.0 \pm 1.8$	$85.6 \pm 5.8$	$76.0 \pm 3.2$	$64.2 \pm 4.3$	$80.2 \pm 2.0$
	MAX-MLP	$73.2 \pm 5.8$	$51.1 \pm 3.6$	–	–	–	$84.0 \pm 6.1$	$76.0 \pm 3.2$	$64.6 \pm 10.2$	$77.8 \pm 1.3$
	MAX-1-LAYER	$72.3 \pm 5.3$	$50.9 \pm 2.2$	–	–	–	$85.1 \pm 7.6$	$75.9 \pm 3.2$	$63.9 \pm 7.7$	$77.7 \pm 1.5$

Classification accuracies

# Conclusion

- Develop theoretical foundations for reasoning about expressive power of GNNs
- Prove tight bounds on the representational capacity of popular GNN variants
- Design a provably most powerful GNN under the message passing framework

Q & A