

EECE 571F: Deep Learning with Structures

Lecture 9: Expressiveness & Generalization of Graph Neural Networks

Renjie Liao

University of British Columbia

Winter, Term 2, 2021/22

Theoretical Aspects of GNNs

- Expressiveness / Capacity

Theoretical Aspects of GNNs

- Expressiveness / Capacity

What functions on graphs can be represented by GNNs?

Can GNNs distinguish isomorphic graphs?

Are GNNs universal approximators?

Theoretical Aspects of GNNs

- Expressiveness / Capacity

What functions on graphs can be represented by GNNs?

Can GNNs distinguish isomorphic graphs?

Are GNNs universal approximators?

Neural Networks, Vol. 2, pp. 359–366, 1989
Printed in the USA. All rights reserved.

0893-6080/89 \$3.00 + .00
Copyright © 1989 Pergamon Press plc

ORIGINAL CONTRIBUTION

Multilayer Feedforward Networks are Universal Approximators

KURT HORNIK

Technische Universität Wien

MAXWELL STINCHCOMBE AND HALBERT WHITE

University of California, San Diego

(Received 16 September 1988; revised and accepted 9 March 1989)

Abstract—*This paper rigorously establishes that standard multilayer feedforward networks with as few as one hidden layer using arbitrary squashing functions are capable of approximating any Borel measurable function from one finite dimensional space to another to any desired degree of accuracy, provided sufficiently many hidden units are available. In this sense, multilayer feedforward networks are a class of universal approximators.*

Theoretical Aspects of GNNs

- Expressiveness / Capacity **Model/Hypothesis Class**

What functions on graphs can be represented by GNNs?

Can GNNs distinguish isomorphic graphs?

Are GNNs universal approximators?

Neural Networks, Vol. 2, pp. 359–366, 1989
Printed in the USA. All rights reserved.

0893-6080/89 \$3.00 + .00
Copyright © 1989 Pergamon Press plc

ORIGINAL CONTRIBUTION

Multilayer Feedforward Networks are Universal Approximators

KURT HORNIK

Technische Universität Wien

MAXWELL STINCHCOMBE AND HALBERT WHITE

University of California, San Diego

(Received 16 September 1988; revised and accepted 9 March 1989)

Abstract—*This paper rigorously establishes that standard multilayer feedforward networks with as few as one hidden layer using arbitrary squashing functions are capable of approximating any Borel measurable function from one finite dimensional space to another to any desired degree of accuracy, provided sufficiently many hidden units are available. In this sense, multilayer feedforward networks are a class of universal approximators.*

Theoretical Aspects of GNNs

- Expressiveness / Capacity **Model/Hypothesis Class**

What functions on graphs can be represented by GNNs?

Can GNNs distinguish isomorphic graphs?

Are GNNs universal approximators?

- Generalization

Theoretical Aspects of GNNs

- Expressiveness / Capacity **Model/Hypothesis Class**

What functions on graphs can be represented by GNNs?

Can GNNs distinguish isomorphic graphs?

Are GNNs universal approximators?

- Generalization

How well do GNNs generalize to unseen graphs?

What are factors that affect the generalization of GNNs?

Theoretical Aspects of GNNs

- Expressiveness / Capacity **Model/Hypothesis Class**

What functions on graphs can be represented by GNNs?

Can GNNs distinguish isomorphic graphs?

Are GNNs universal approximators?

- Generalization **Model/Hypothesis Class + Learning Process**

How well do GNNs generalize to unseen graphs?

What are factors that affect the generalization of GNNs?

Theoretical Aspects of GNNs

- **Expressiveness / Capacity**

What functions on graphs can be represented by GNNs?

Can GNNs distinguish isomorphic graphs?

Are GNNs universal approximators?

- **Generalization**

How well do GNNs generalize to unseen graphs?

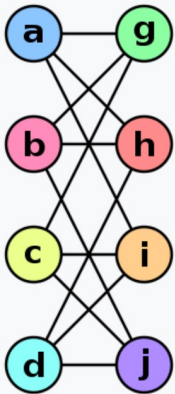
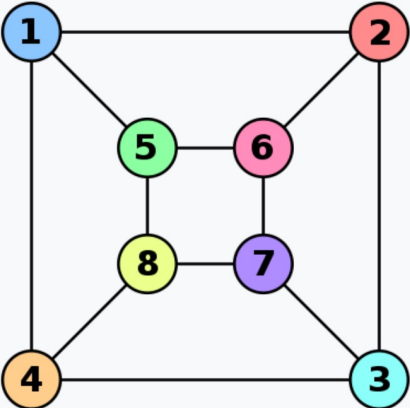
What are factors that affect the generalization of GNNs?

Graph Isomorphism Problem

Given two graphs, are they isomorphic?

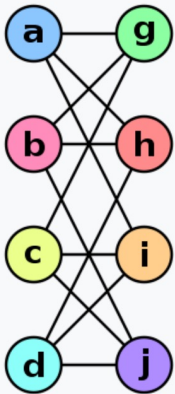
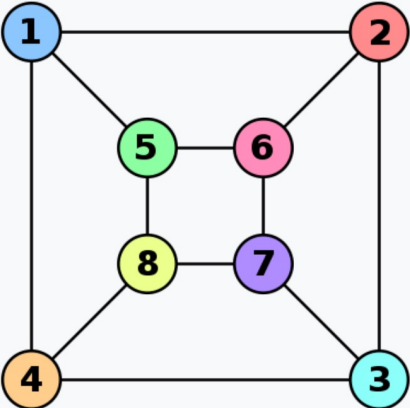
Graph Isomorphism Problem

Given two graphs, are they isomorphic?

Graph G	Graph H	An isomorphism between G and H
		$f(a) = 1$ $f(b) = 6$ $f(c) = 8$ $f(d) = 3$ $f(g) = 5$ $f(h) = 2$ $f(i) = 4$ $f(j) = 7$

Graph Isomorphism Problem

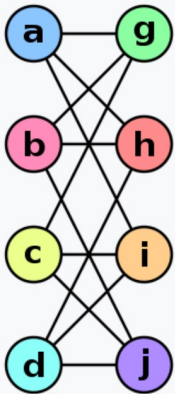
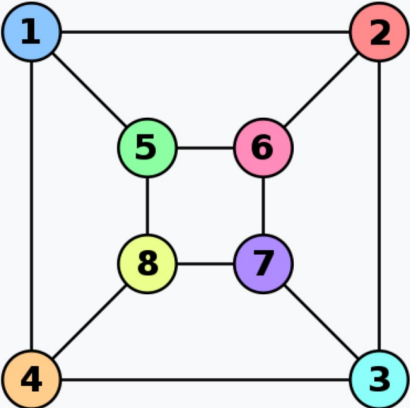
Given two graphs, are they isomorphic?

Graph G	Graph H	An isomorphism between G and H
		$f(a) = 1$ $f(b) = 6$ $f(c) = 8$ $f(d) = 3$ $f(g) = 5$ $f(h) = 2$ $f(i) = 4$ $f(j) = 7$

Given two adjacency matrices A_1, A_2 , is there a permutation matrix P s.t. $PA_1P^T = A_2$?

Graph Isomorphism Problem

Given two graphs, are they isomorphic?

Graph G	Graph H	An isomorphism between G and H
		$f(a) = 1$ $f(b) = 6$ $f(c) = 8$ $f(d) = 3$ $f(g) = 5$ $f(h) = 2$ $f(i) = 4$ $f(j) = 7$

It is unknown if this problem is in P or NP-Complete!

Given two adjacency matrices A_1, A_2 , is there a permutation matrix P s.t. $PA_1P^T = A_2$?

The Weisfeiler-Lehman Isomorphism Test

Algorithm 1: 1-WL Algorithm [1]

Input: Initial node label $(h_1^{(0)}, h_2^{(0)}, \dots, h_N^{(0)})$

$t \leftarrow 0$;

repeat

for $v_i \in \mathcal{V}$ **do**

$h_i^{(t+1)} \leftarrow \text{hash}(\{h_j^{(t)} \mid j \in \mathcal{N}_i\})$;

$t \leftarrow t + 1$;

until *stable node partition based on labels are reached or $t = N$* ;

Output: Final node label $(h_1^{(T)}, h_2^{(T)}, \dots, h_N^{(T)})$

The Weisfeiler-Lehman Isomorphism Test

Algorithm 1: 1-WL Algorithm [1]

Input: Initial node label $(h_1^{(0)}, h_2^{(0)}, \dots, h_N^{(0)})$

$t \leftarrow 0$;

repeat

for $v_i \in \mathcal{V}$ **do**

$h_i^{(t+1)} \leftarrow \text{hash}(\{h_j^{(t)} \mid j \in \mathcal{N}_i\})$;

$t \leftarrow t + 1$;

until *stable node partition based on labels are reached or* $t = N$;

Output: Final node label $(h_1^{(T)}, h_2^{(T)}, \dots, h_N^{(T)})$

Injective hash function:

Mapping distinct multi-sets to distinct labels

The Weisfeiler-Lehman Isomorphism Test

Algorithm 1: 1-WL Algorithm [1]

Input: Initial node label $(h_1^{(0)}, h_2^{(0)}, \dots, h_N^{(0)})$

$t \leftarrow 0$;

repeat

for $v_i \in \mathcal{V}$ **do**

$h_i^{(t+1)} \leftarrow \text{hash}(\{h_j^{(t)} \mid j \in \mathcal{N}_i\})$;

$t \leftarrow t + 1$;

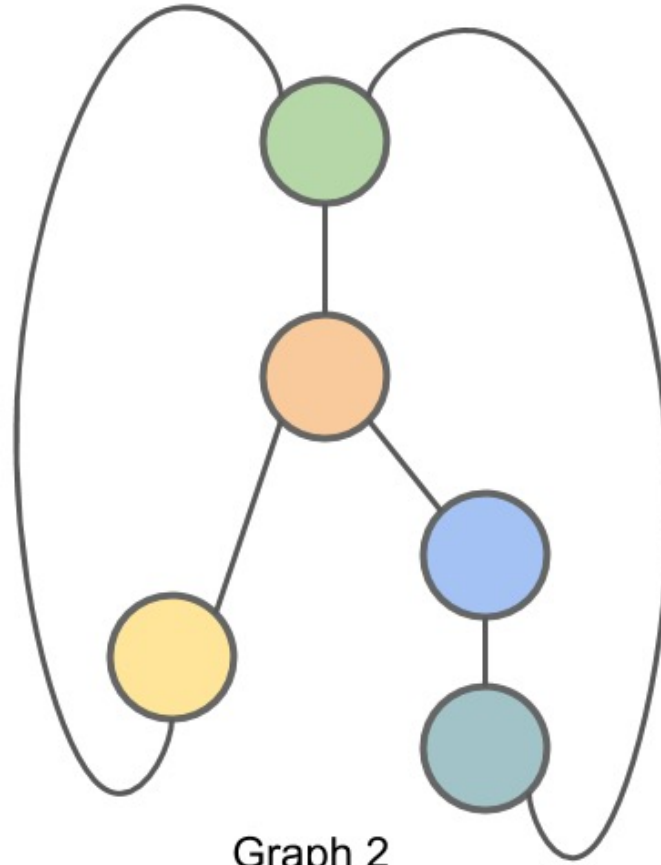
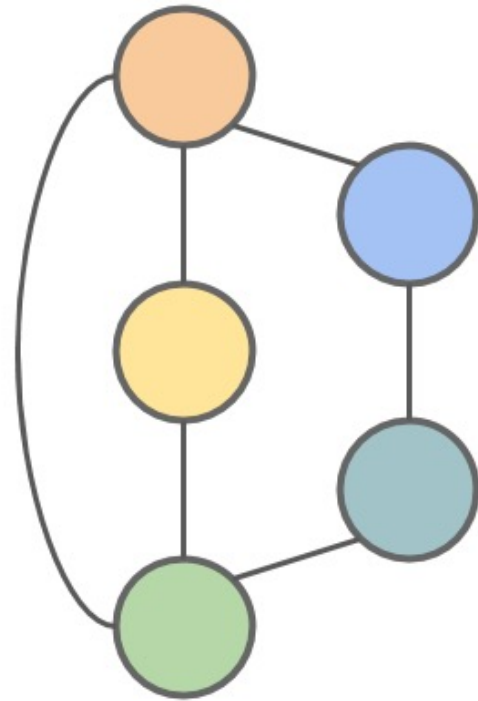
until *stable node partition based on labels are reached or $t = N$* ;

Output: Final node label $(h_1^{(T)}, h_2^{(T)}, \dots, h_N^{(T)})$

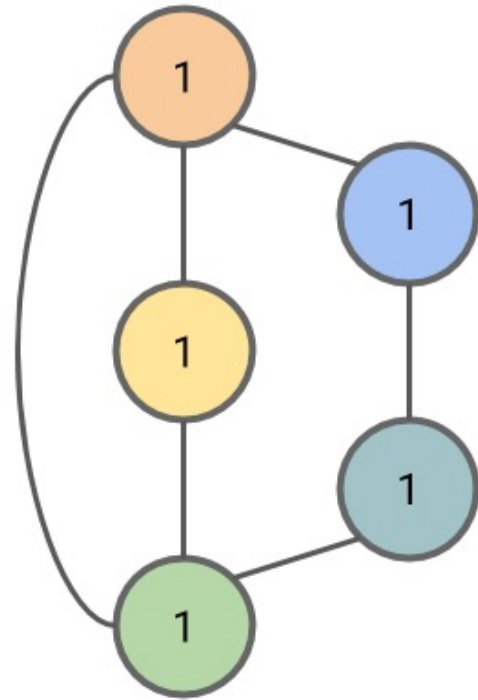
If the partition of nodes by labels are different, then two graphs are non-isomorphic!

If the partition of nodes by labels are same, we can not decide!

The Weisfeiler-Lehman Isomorphism Test

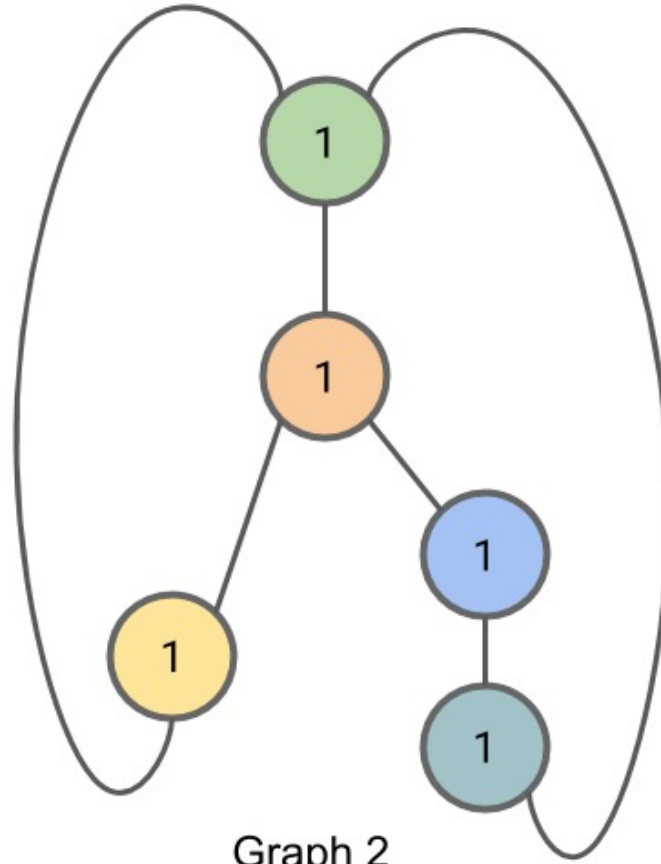


The Weisfeiler-Lehman Isomorphism Test



Graph 1

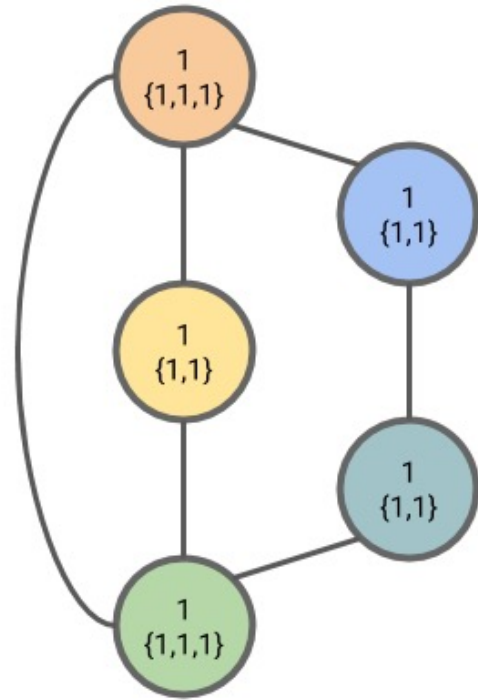
$\{1,1,1,1,1\}$



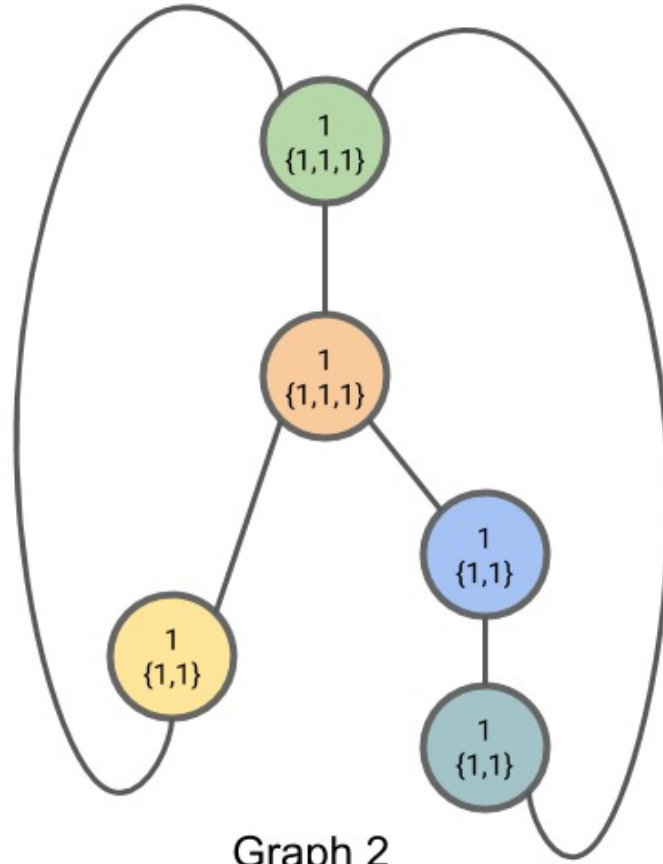
Graph 2

$\{1,1,1,1,1\}$

The Weisfeiler-Lehman Isomorphism Test



Graph 1



Graph 2

Hash Table:

$\{1,1\} \rightarrow 2$

$\{1,1,1\} \rightarrow 3$

$\{2,3\} \rightarrow 4$

$\{3,3\} \rightarrow 5$

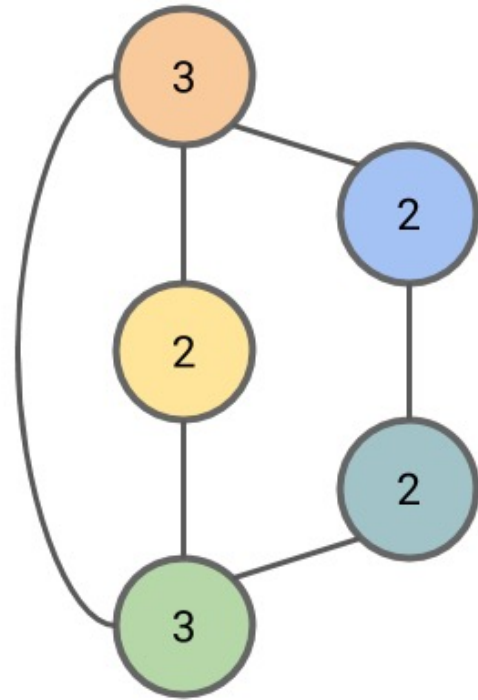
$\{2,2,3\} \rightarrow 6$

$\{4,6\} \rightarrow 7$

$\{6,6\} \rightarrow 8$

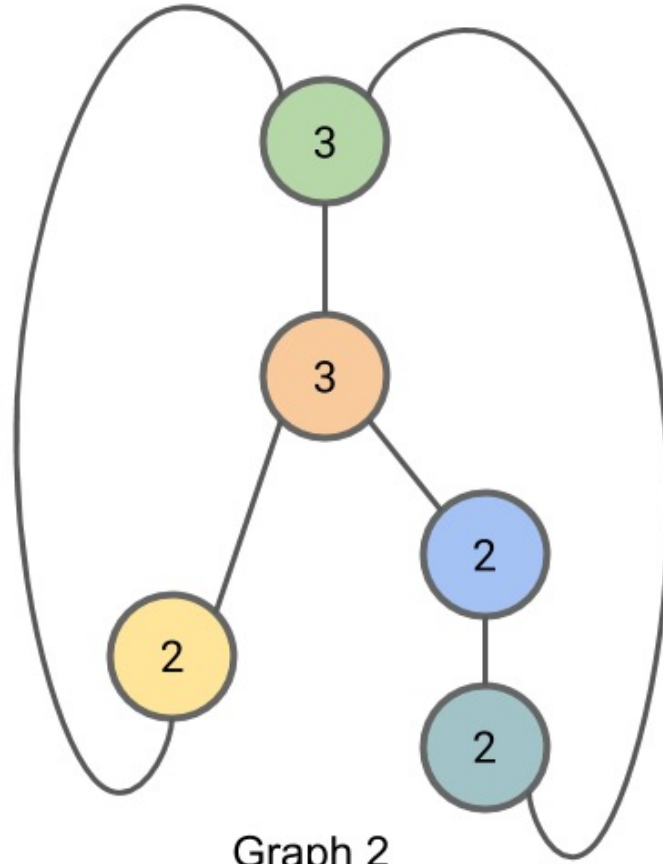
$\{4,5,6\} \rightarrow 9$

The Weisfeiler-Lehman Isomorphism Test



Graph 1

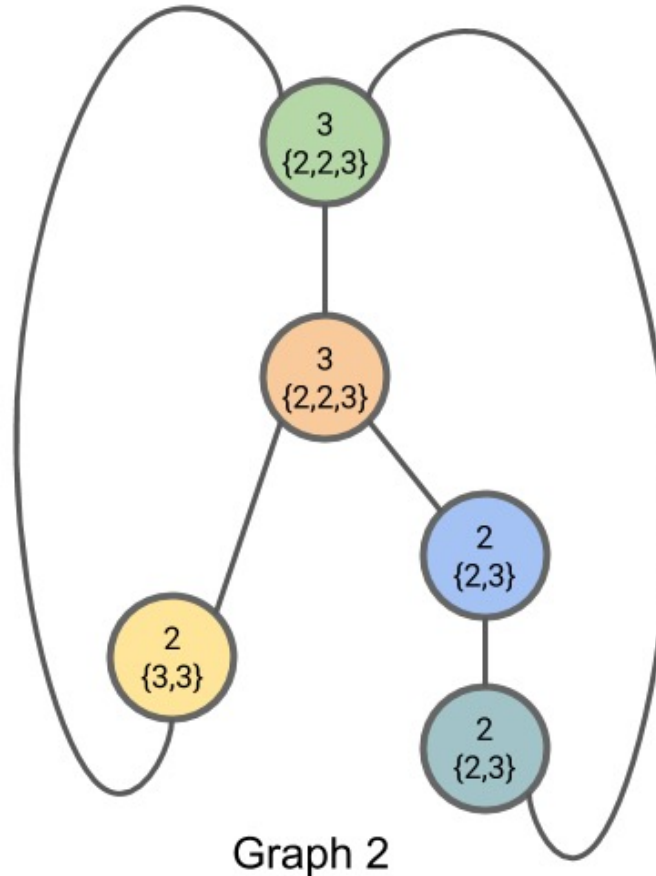
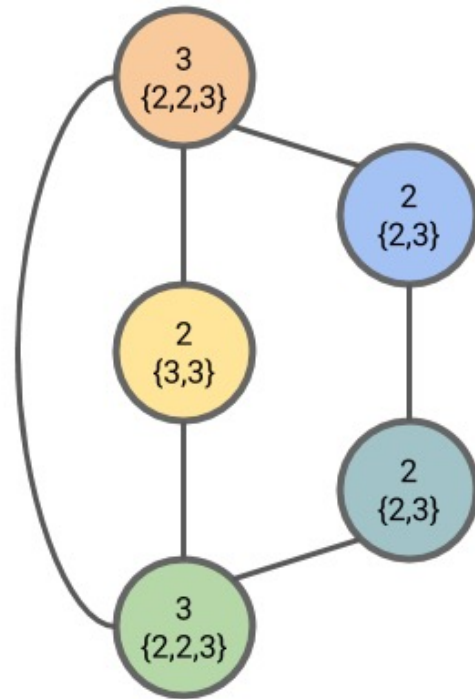
$\{2,2,2\}, \{3,3\}$



Graph 2

$\{2,2,2\}, \{3,3\}$

The Weisfeiler-Lehman Isomorphism Test



Hash Table:

{1,1} -> 2

{1,1,1} -> 3

{2,3} -> 4

{3,3} -> 5

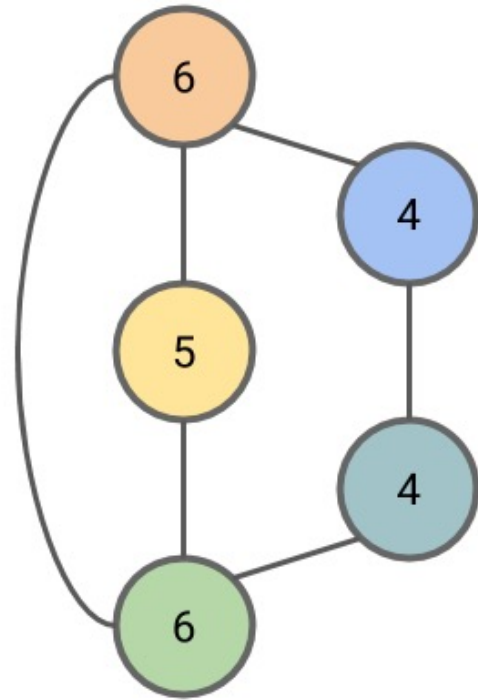
{2,2,3} -> 6

{4,6} -> 7

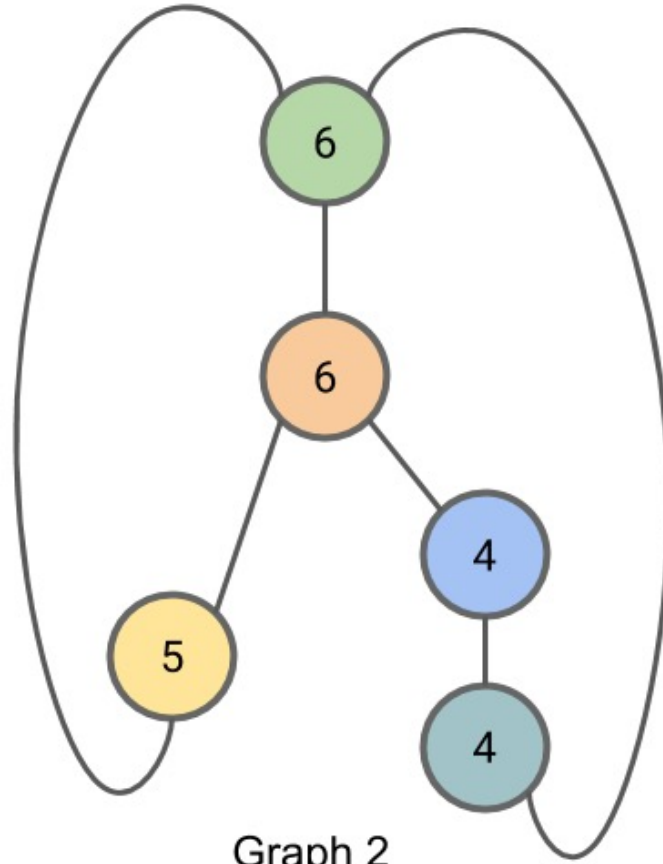
{6,6} -> 8

{4,5,6} -> 9

The Weisfeiler-Lehman Isomorphism Test

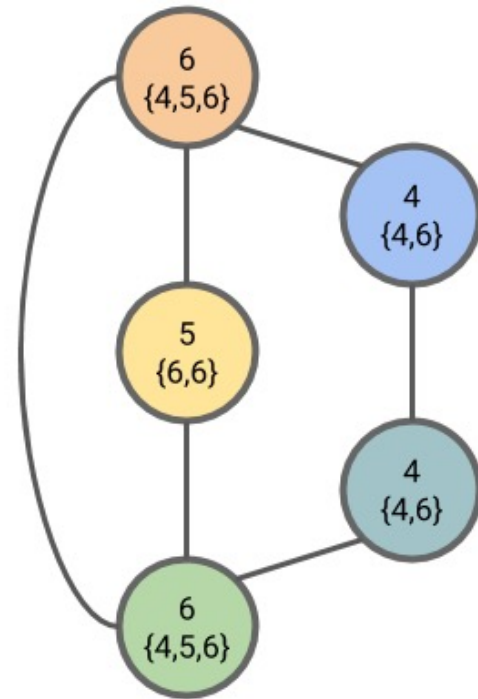


$\{4,4\}, \{5\}, \{6,6\}$

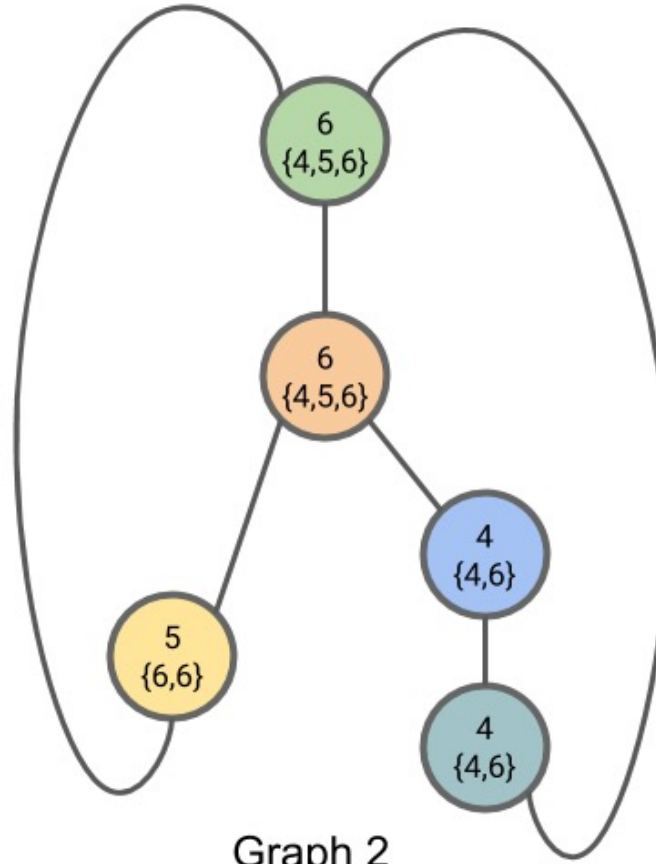


$\{4,4\}, \{5\}, \{6,6\}$

The Weisfeiler-Lehman Isomorphism Test



Graph 1

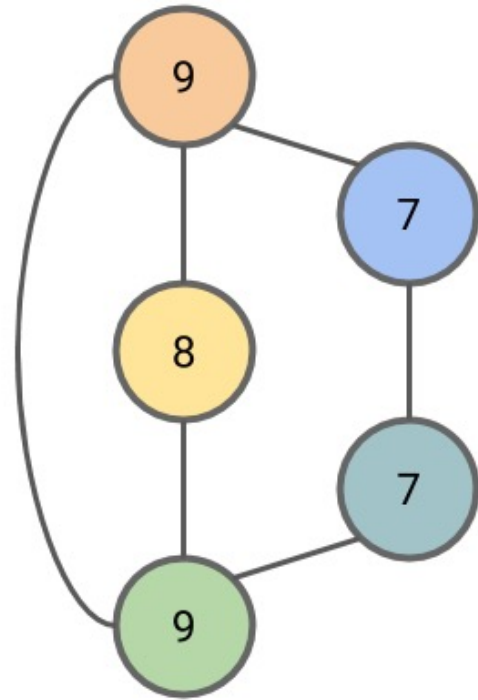


Graph 2

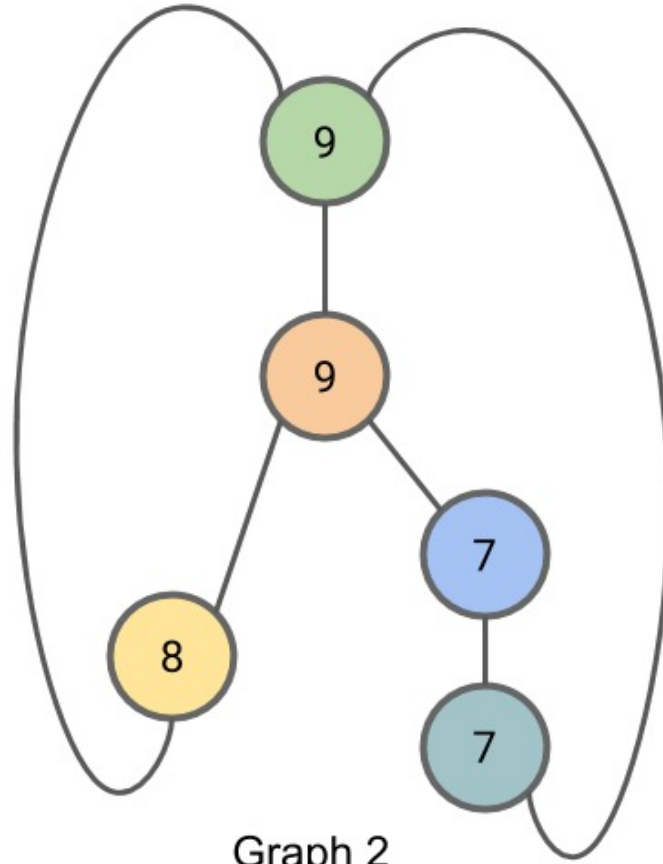
Hash Table:

$\{1,1\} \rightarrow 2$
 $\{1,1,1\} \rightarrow 3$
 $\{2,3\} \rightarrow 4$
 $\{3,3\} \rightarrow 5$
 $\{2,2,3\} \rightarrow 6$
 $\{4,6\} \rightarrow 7$
 $\{6,6\} \rightarrow 8$
 $\{4,5,6\} \rightarrow 9$

The Weisfeiler-Lehman Isomorphism Test



$\{7,7\}, \{8\}, \{9,9\}$



$\{7,7\}, \{8\}, \{9,9\}$

A Recap on GNNs

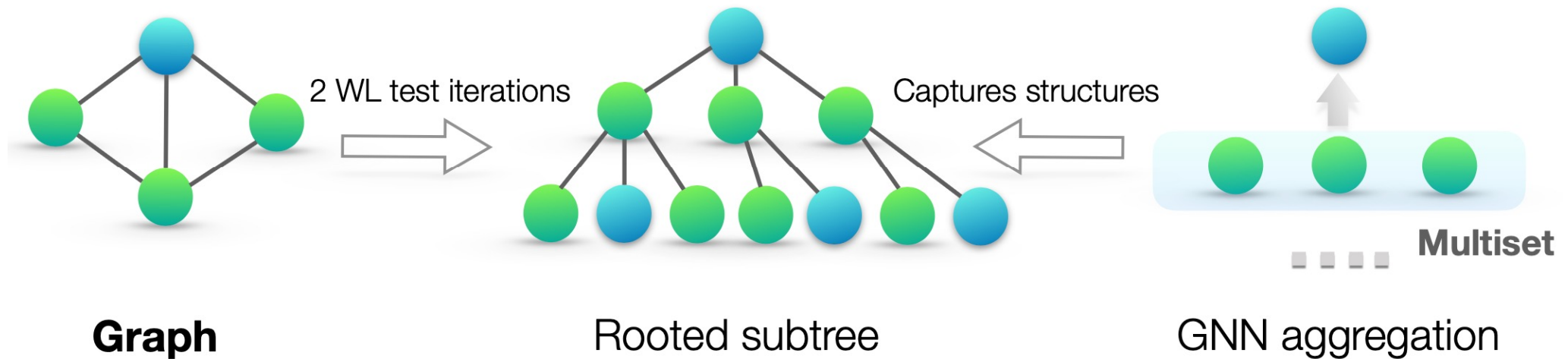
$$a_v^{(k)} = \text{AGGREGATE}^{(k)} \left(\left\{ h_u^{(k-1)} : u \in \mathcal{N}(v) \right\} \right)$$

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

A Recap on GNNs

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

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



GNNs Are as Powerful as 1-WL

Lemma 2. *Let G_1 and G_2 be any two non-isomorphic graphs. If a graph neural network $\mathcal{A} : \mathcal{G} \rightarrow \mathbb{R}^d$ 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 as Powerful as 1-WL

Lemma 2. *Let G_1 and G_2 be any two non-isomorphic graphs. If a graph neural network $\mathcal{A} : \mathcal{G} \rightarrow \mathbb{R}^d$ 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.*

Proof Sketch: Proof by contradiction.

Suppose the GNN has $\mathcal{A}(G_1) \neq \mathcal{A}(G_2)$, but WL test outputs the same node partition based on labels

By induction, prove that there exists a valid mapping from node label (WL) to node feature (GNN)

Readout of GNN is permutation invariant \Rightarrow Same node partition generates same graph representation

Therefore, we have $\mathcal{A}(G_1) = \mathcal{A}(G_2)$

GNNs Are as Powerful as 1-WL

Theorem 3. *Let $\mathcal{A} : \mathcal{G} \rightarrow \mathbb{R}^d$ be a GNN. With a sufficient number of GNN layers, \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 ϕ are injective.

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

GNNs Are as Powerful as 1-WL

Theorem 3. *Let $\mathcal{A} : \mathcal{G} \rightarrow \mathbb{R}^d$ be a GNN. With a sufficient number of GNN layers, \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 ϕ are injective.

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

Proof Sketch:

By induction, prove that there exists an injective mapping from node label (WL) to node feature (GNN)

Since readout is injective, different multi-sets will be mapped to different graph embeddings

Graph Isomorphism Networks (GINs)

Lemma 5. *Assume \mathcal{X} is countable. There exists a function $f : \mathcal{X} \rightarrow \mathbb{R}^n$ so that $h(X) = \sum_{x \in X} f(x)$ is unique for each multiset $X \subset \mathcal{X}$ of bounded size. Moreover, any multiset function g can be decomposed as $g(X) = \phi(\sum_{x \in X} f(x))$ for some function ϕ .*

Graph Isomorphism Networks (GINs)

Lemma 5. *Assume \mathcal{X} is countable. There exists a function $f : \mathcal{X} \rightarrow \mathbb{R}^n$ so that $h(X) = \sum_{x \in X} f(x)$ is unique for each multiset $X \subset \mathcal{X}$ of bounded size. Moreover, any multiset function g can be decomposed as $g(X) = \phi(\sum_{x \in X} f(x))$ for some function ϕ .*

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 ϵ , 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 multiset of bounded size. Moreover, any function g over such pairs can be decomposed as $g(c, X) = \varphi((1 + \epsilon) \cdot f(c) + \sum_{x \in X} f(x))$ for some function φ .*

Graph Isomorphism Networks (GINs)

Lemma 5. *Assume \mathcal{X} is countable. There exists a function $f : \mathcal{X} \rightarrow \mathbb{R}^n$ so that $h(X) = \sum_{x \in X} f(x)$ is unique for each multiset $X \subset \mathcal{X}$ of bounded size. Moreover, any multiset function g can be decomposed as $g(X) = \phi\left(\sum_{x \in X} f(x)\right)$ for some function ϕ .*

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 ϵ , 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 multiset of bounded size. Moreover, any function g over such pairs can be decomposed as $g(c, X) = \varphi\left((1 + \epsilon) \cdot f(c) + \sum_{x \in X} f(x)\right)$ for some function φ .*

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

Other Results & Open Questions?

- GNNs are universal approximators for functions on graphs [5,6]
- GNNs are Turing universal [7]
- Expressiveness of GNNs using *communication capacity* [8]
- Expressiveness of GNNs in terms of *counting substructures* [9]
- Can we design GNNs that go beyond 1-WL? [3,17,18,19]

Theoretical Aspects of GNNs

- Expressiveness / Capacity

What functions on graphs can be represented by GNNs?

Can GNNs distinguish isomorphic graphs?

Are GNNs universal approximators?

- **Generalization**

How well do GNNs generalize to unseen graphs?

What are factors that affect the generalization of GNNs?

Statistical Learning Theory

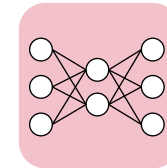
i.i.d. data



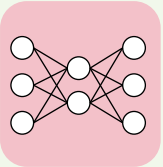
Learning Algorithm



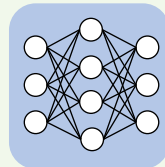
best model



hypothesis class \mathcal{H}



.....

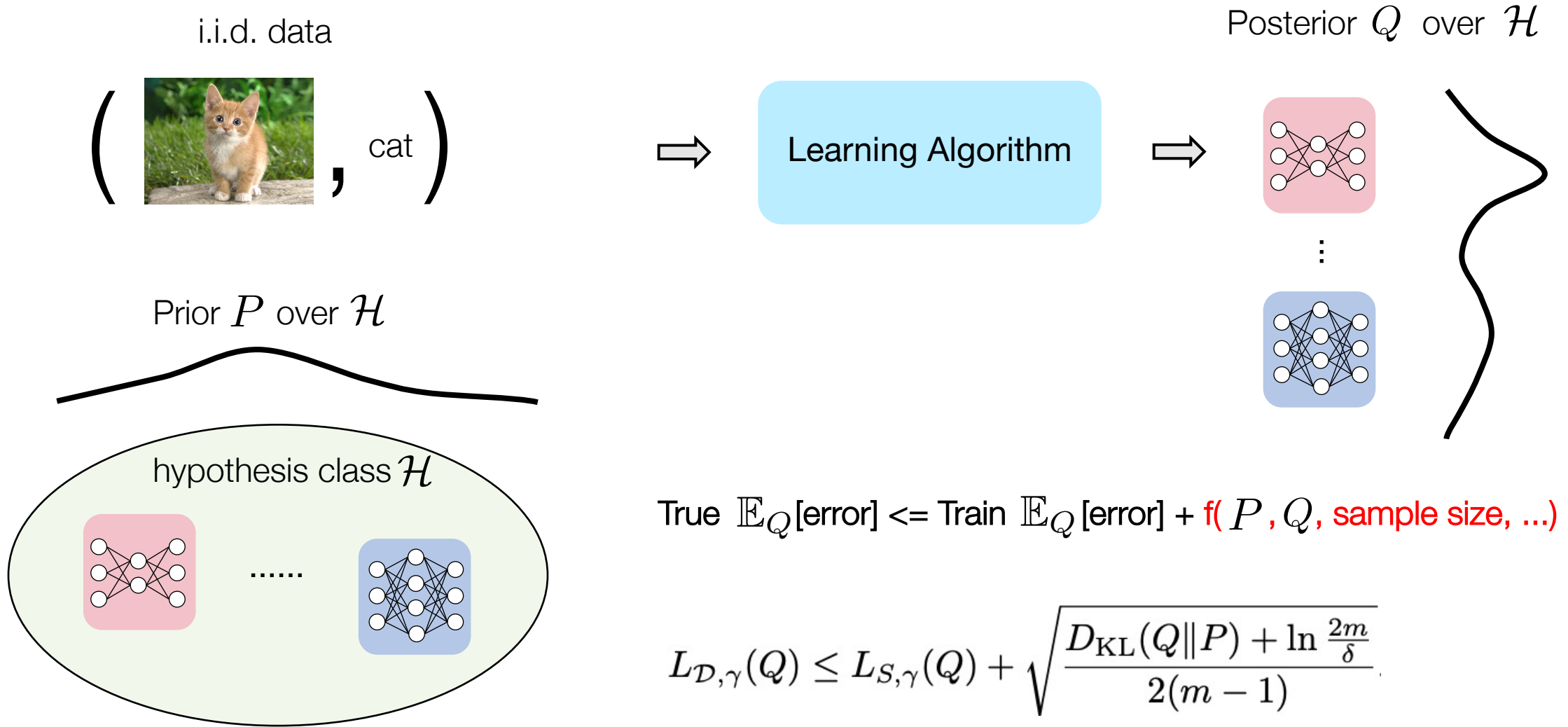


(Uniform Convergence) Generalization Bound:

For any hypothesis in \mathcal{H} , any data distribution, any $\delta \in (0, 1)$,
with probability $1 - \delta$,

True error \leq Train error + **f(complexity of \mathcal{H} , sample size, ...)**

PAC-Bayes Theory



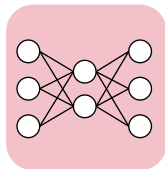
PAC-Bayes Theory

Lemma 2.2. (Neyshabur et al., 2017)⁴ Let $f_w(x) : \mathcal{X} \rightarrow \mathbb{R}^K$ be any model with parameters w , and let P be any distribution on the parameters that is independent of the training data. For any w , we construct a posterior $Q(w + u)$ by adding any random perturbation u to w , s.t., $\mathbb{P}(\max_{x \in \mathcal{X}} |f_{w+u}(x) - f_w(x)|_\infty < \frac{\gamma}{4}) > \frac{1}{2}$. Then, for any $\gamma, \delta > 0$, with probability at least $1 - \delta$ over an i.i.d. size- m training set S according to \mathcal{D} , for any w , we have:

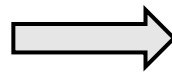
$$L_{\mathcal{D},0}(f_w) \leq L_{S,\gamma}(f_w) + \sqrt{\frac{2D_{\text{KL}}(Q(w+u)\|P) + \log \frac{8m}{\delta}}{2(m-1)}}.$$

True error \leq Train error + **f(P , Q , sample size, ...)**

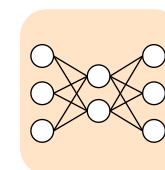
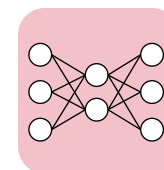
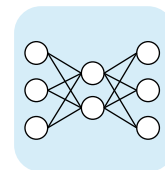
Learned Model



Perturb Weights

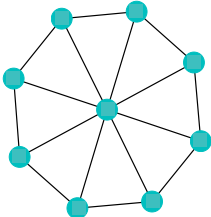


Posterior Q over \mathcal{H}

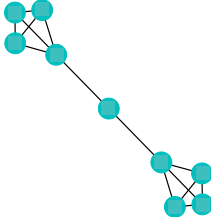


Problems & Assumptions

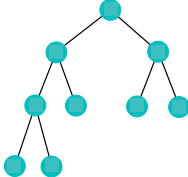
Graph Classification



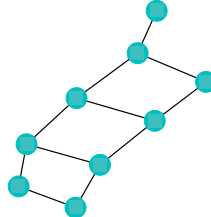
Wheel



Barbell



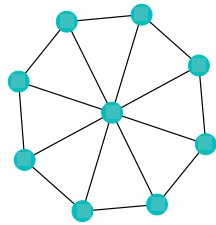
Binary Tree



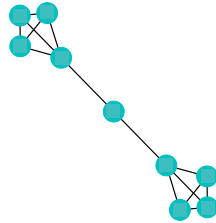
Ladder

Problems & Assumptions

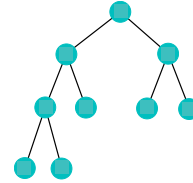
Graph Classification



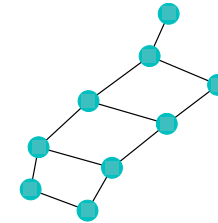
Wheel



Barbell



Binary Tree



Ladder

- A1 Data, *i.e.*, triplets (A, X, y) , are i.i.d. samples drawn from some unknown distribution \mathcal{D} .
- A2 The maximum hidden dimension across all layers is h .
- A3 Node feature of any graph is contained in a ℓ_2 -ball with radius B . Specifically, we have $\forall i \in \mathbb{N}_n^+$, the i -th node feature $X[i, :] \in \mathcal{X}_{B, h_0} = \{x \in \mathbb{R}^{h_0} \mid \sum_{j=1}^{h_0} x_j^2 \leq B^2\}$.
- A4 We only consider simple graphs (*i.e.*, undirected, no loops¹, and no multi-edges) with maximum node degree as $d - 1$.

Models

Graph Convolutional Networks (GCNs) [12]

$$H_k = \sigma_k \left(\tilde{L} H_{k-1} W_k \right) \quad (k\text{-th Graph Convolution Layer})$$

$$H_l = \frac{1}{n} \mathbf{1}_n H_{l-1} W_l \quad (\text{Readout Layer}),$$

Models

Graph Convolutional Networks (GCNs) [12]

$$H_k = \sigma_k \left(\tilde{L} H_{k-1} W_k \right) \quad (k\text{-th Graph Convolution Layer})$$

$$H_l = \frac{1}{n} \mathbf{1}_n H_{l-1} W_l \quad (\text{Readout Layer}),$$

Message Passing Graph Neural Networks (MPGNNs) [13]

$$M_k = g(C_{\text{out}}^\top H_{k-1}) \quad (k\text{-th step Message Computation})$$

$$\bar{M}_k = C_{\text{in}} M_k \quad (k\text{-th step Message Aggregation})$$

$$H_k = \phi \left(X W_1 + \rho \left(\bar{M}_k \right) W_2 \right) \quad (k\text{-th step Node State Update})$$

$$H_l = \frac{1}{n} \mathbf{1}_n H_{l-1} W_l \quad (\text{Readout Layer}),$$

Perturbation Bounds

Lemma 3.1. (*GCN Perturbation Bound*) For any $B > 0, l > 1$, let $f_w \in \mathcal{H} : \mathcal{X} \times \mathcal{G} \rightarrow \mathbb{R}^K$ be a l -layer GCN. Then for any w , and $x \in \mathcal{X}_{B, h_0}$, and any perturbation $u = \text{vec}(\{U_i\}_{i=1}^l)$ such that $\forall i \in \mathbb{N}_l^+, \|U_i\|_2 \leq \frac{1}{l} \|W_i\|_2$, the change in the output of GCN is bounded as,

$$|f_{w+u}(X, A) - f_w(X, A)|_2 \leq eBd^{\frac{l-1}{2}} \left(\prod_{i=1}^l \|W_i\|_2 \right) \sum_{k=1}^l \frac{\|U_k\|_2}{\|W_k\|_2}.$$

Perturbation Bounds

Lemma 3.1. (*GCN Perturbation Bound*) For any $B > 0, l > 1$, let $f_w \in \mathcal{H} : \mathcal{X} \times \mathcal{G} \rightarrow \mathbb{R}^K$ be a l -layer GCN. Then for any w , and $x \in \mathcal{X}_{B, h_0}$, and any perturbation $u = \text{vec}(\{U_i\}_{i=1}^l)$ such that $\forall i \in \mathbb{N}_l^+, \|U_i\|_2 \leq \frac{1}{l} \|W_i\|_2$, the change in the output of GCN is bounded as,

$$|f_{w+u}(X, A) - f_w(X, A)|_2 \leq eBd^{\frac{l-1}{2}} \left(\prod_{i=1}^l \|W_i\|_2 \right) \sum_{k=1}^l \frac{\|U_k\|_2}{\|W_k\|_2}.$$

Lemma 3.3. (*MPGNN Perturbation Bound*) For any $B > 0, l > 1$, let $f_w \in \mathcal{H} : \mathcal{X} \times \mathcal{G} \rightarrow \mathbb{R}^K$ be a l -step MPGNN. Then for any w , and $x \in \mathcal{X}_{B, h_0}$, and any perturbation $u = \text{vec}(\{U_1, U_2, U_l\})$ such that $\eta = \max\left(\frac{\|U_1\|_2}{\|W_1\|_2}, \frac{\|U_2\|_2}{\|W_2\|_2}, \frac{\|U_l\|_2}{\|W_l\|_2}\right) \leq \frac{1}{l}$, the change in the output of MPGNN is bounded as,

$$|f_{w+u}(X, A) - f_w(X, A)|_2 \leq eBl\eta \|W_1\|_2 \|W_l\|_2 C_\phi \frac{(d\mathcal{C})^{l-1} - 1}{d\mathcal{C} - 1},$$

where $\mathcal{C} = C_\phi C_\rho C_g \|W_2\|_2$.

Generalization Bounds for GCNs

Theorem 3.2. (*GCN Generalization Bound*) For any $B > 0, l > 1$, let $f_w \in \mathcal{H} : \mathcal{X} \times \mathcal{G} \rightarrow \mathbb{R}^K$ be a l layer GCN. Then for any $\delta, \gamma > 0$, with probability at least $1 - \delta$ over the choice of an i.i.d. size- m training set S according to \mathcal{D} , for any w , we have,

$$L_{\mathcal{D},0}(f_w) \leq L_{S,\gamma}(f_w) + \mathcal{O} \left(\sqrt{\frac{B^2 d^{l-1} l^2 h \log(lh) \prod_{i=1}^l \|W_i\|_2^2 \sum_{i=1}^l (\|W_i\|_F^2 / \|W_i\|_2^2) + \log \frac{ml}{\delta}}{\gamma^2 m}} \right).$$

Generalization Bounds for GCNs

Theorem 3.2. (GCN Generalization Bound) For any $B > 0, l > 1$, let $f_w \in \mathcal{H} : \mathcal{X} \times \mathcal{G} \rightarrow \mathbb{R}^K$ be a l layer GCN. Then for any $\delta, \gamma > 0$, with probability at least $1 - \delta$ over the choice of an i.i.d. size- m training set S according to \mathcal{D} , for any w , we have,

$$L_{\mathcal{D},0}(f_w) \leq L_{S,\gamma}(f_w) + \mathcal{O} \left(\sqrt{\frac{B^2 d^{l-1} l^2 h \log(lh) \prod_{i=1}^l \|W_i\|_2^2 \sum_{i=1}^l (\|W_i\|_F^2 / \|W_i\|_2^2) + \log \frac{ml}{\delta}}{\gamma^2 m}} \right).$$

Proof Sketch:

Perturbation bound + Measure Concentration Inequalities \Rightarrow A model with certain learned weights could generalize

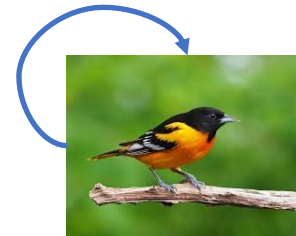
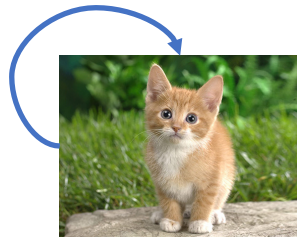
Similar bounds hold for a finite covering set of all possible weights and then use union bound to derive the final result

Generalization Bounds for GCNs

MLPs/CNNs are special GCNs if we construct the following graph:

- Each image is a node
- No edges exist except self-loops

(Convolution is matrix multiplication as well)



Generalization Bounds for GCNs

Theorem 3.2. (GCN Generalization Bound) For any $B > 0, l > 1$, let $f_w \in \mathcal{H} : \mathcal{X} \times \mathcal{G} \rightarrow \mathbb{R}^K$ be a l layer GCN. Then for any $\delta, \gamma > 0$, with probability at least $1 - \delta$ over the choice of an i.i.d. size- m training set S according to \mathcal{D} , for any w , we have,

$$L_{\mathcal{D},0}(f_w) \leq L_{S,\gamma}(f_w) + \mathcal{O} \left(\sqrt{\frac{B^2 d^{l-1} l^2 h \log(lh) \prod_{i=1}^l \|W_i\|_2^2 \sum_{i=1}^l (\|W_i\|_F^2 / \|W_i\|_2^2) + \log \frac{ml}{\delta}}{\gamma^2 m}} \right).$$

PAC-Bayes bounds for MLPs/CNNs (w/ ReLU)

$$L_{D,0}(f_w) \leq L_{S,\gamma}(f_w) + \mathcal{O} \left(\sqrt{\frac{B^2 l^2 h \log(lh) \prod_{i=1}^l \|W_i\|_2^2 \sum_{i=1}^l (\|W_i\|_F^2 / \|W_i\|_2^2) + \log \frac{ml}{\delta}}{\gamma^2 m}} \right).$$

Generalization Bounds for GCNs

Theorem 3.2. (*GCN Generalization Bound*) For any $B > 0, l > 1$, let $f_w \in \mathcal{H} : \mathcal{X} \times \mathcal{G} \rightarrow \mathbb{R}^K$ be a l layer GCN. Then for any $\delta, \gamma > 0$, with probability at least $1 - \delta$ over the choice of an i.i.d. size- m training set S according to \mathcal{D} , for any w , we have,

$$L_{\mathcal{D},0}(f_w) \leq L_{S,\gamma}(f_w) + \mathcal{O} \left(\sqrt{\frac{B^2 d^{l-1} l^2 h \log(lh) \prod_{i=1}^l \|W_i\|_2^2 \sum_{i=1}^l (\|W_i\|_F^2 / \|W_i\|_2^2) + \log \frac{ml}{\delta}}{\gamma^2 m}} \right).$$

PAC-Bayes bounds for MLPs/CNNs (w/ ReLU)

$$L_{D,0}(f_w) \leq L_{S,\gamma}(f_w) + \mathcal{O} \left(\sqrt{\frac{B^2 l^2 h \log(lh) \prod_{i=1}^l \|W_i\|_2^2 \sum_{i=1}^l (\|W_i\|_F^2 / \|W_i\|_2^2) + \log \frac{ml}{\delta}}{\gamma^2 m}} \right).$$

Generalization Bounds for Message Passing GNNs

Theorem 3.4. (*MPGNN Generalization Bound*) For any $B > 0, l > 1$, let $f_w \in \mathcal{H} : \mathcal{X} \times \mathcal{G} \rightarrow \mathbb{R}^K$ be a l -step MPGNN. Then for any $\delta, \gamma > 0$, with probability at least $1 - \delta$ over the choice of an i.i.d. size- m training set S according to \mathcal{D} , for any w , we have,

$$L_{\mathcal{D},0}(f_w) \leq L_{S,\gamma}(f_w) + \mathcal{O} \left(\sqrt{\frac{B^2 (\max(\zeta^{-(l+1)}, (\lambda\xi)^{(l+1)/l}))^2 l^2 h \log(lh) |w|_2^2 + \log \frac{m(l+1)}{\delta}}{\gamma^2 m}} \right),$$

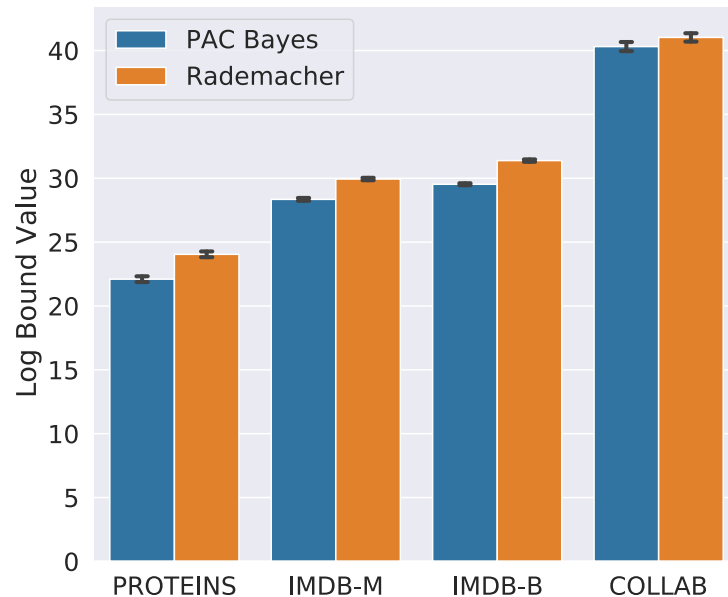
where $\zeta = \min(\|W_1\|_2, \|W_2\|_2, \|W_l\|_2)$, $|w|_2^2 = \|W_1\|_F^2 + \|W_2\|_F^2 + \|W_l\|_F^2$, $\mathcal{C} = C_\phi C_\rho C_g \|W_2\|_2$, $\lambda = \|W_1\|_2 \|W_l\|_2$, and $\xi = C_\phi \frac{(d\mathcal{C})^{l-1} - 1}{d\mathcal{C} - 1}$.

Generalization Bounds for Message Passing GNNs

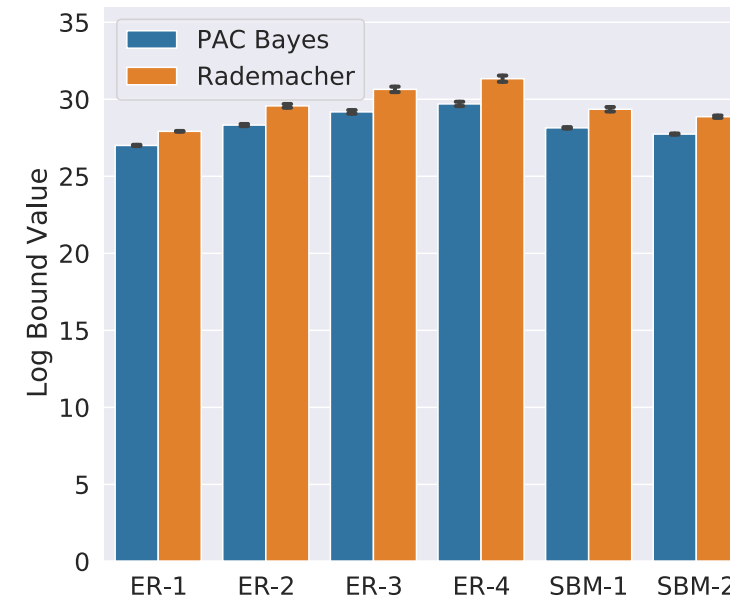
Statistics	Max Node Degree $d - 1$	Max Hidden Dim h	Spectral Norm of Learned Weights
VC-Dimension (Scarselli et al., 2018)	-	$\mathcal{O}(h^4)$	-
Rademacher Complexity (Garg et al., 2020)	$\mathcal{O}(d^{l-1} \sqrt{\log(d^{2l-3})})$	$\mathcal{O}(h\sqrt{\log h})$	$\mathcal{O}(\lambda \mathcal{C} \xi \sqrt{\log(\ W_2\ _2 \lambda \xi^2)})$
Ours	$\mathcal{O}(d^{l-1})$	$\mathcal{O}(\sqrt{h \log h})$	$\mathcal{O}(\lambda^{1+\frac{1}{l}} \xi^{1+\frac{1}{l}} \sqrt{\ W_1\ _F^2 + \ W_2\ _F^2 + \ W_l\ _F^2})$

Generalization Bounds for Message Passing GNNs

Statistics	Max Node Degree $d - 1$	Max Hidden Dim h	Spectral Norm of Learned Weights
VC-Dimension (Scarselli et al., 2018)	-	$\mathcal{O}(h^4)$	-
Rademacher Complexity (Garg et al., 2020)	$\mathcal{O}(d^{l-1} \sqrt{\log(d^{2l-3})})$	$\mathcal{O}(h\sqrt{\log h})$	$\mathcal{O}(\lambda \mathcal{C} \xi \sqrt{\log(\ W_2\ _2 \lambda \xi^2)})$
Ours	$\mathcal{O}(d^{l-1})$	$\mathcal{O}(\sqrt{h \log h})$	$\mathcal{O}(\lambda^{1+\frac{1}{l}} \xi^{1+\frac{1}{l}} \sqrt{\ W_1\ _F^2 + \ W_2\ _F^2 + \ W_l\ _F^2})$



real-world graphs



synthetic graphs

Take Home Messages & Other Results

- Spectral norm of weights and max node degree play key roles in the generalization of GNNs
- Bounds for MLPs/CNNs (w/ ReLU) are special cases of bounds for GCNs, which reconciles with the observation that MLPs/CNNs (w/ ReLU) can be viewed as special GCNs
- PAC-Bayes bounds are empirically tighter than the recent Rademacher complexity based bounds on several synthetic and real-world graph datasets
- Rademacher complexity based generalization bounds [15]
- VC-Dimension based generalization bounds [16]

Open Questions

- Is the maximum node degree the only graph statistic that has an impact on the generalization ability of GNNs?
- Would the analysis still work for other interesting problem setups like out-of-distribution generalizations?
- What is the impact of the optimization algorithms like SGD on the generalization of GNNs?

References

- [1] Weisfeiler, B. and Leman, A., 1968. The reduction of a graph to canonical form and the algebra which appears therein. *NTI, Series*, 2(9), pp.12-16.
- [2] Xu, K., Hu, W., Leskovec, J., & Jegelka, S. (2018). How powerful are graph neural networks?. *arXiv preprint arXiv:1810.00826*.
- [3] Morris, C., Ritzert, M., Fey, M., Hamilton, W. L., Lenssen, J. E., Rattan, G., & Grohe, M. (2019, July). Weisfeiler and leman go neural: Higher-order graph neural networks. In *Proceedings of the AAAI conference on artificial intelligence* (Vol. 33, No. 01, pp. 4602-4609).
- [4] Shervashidze, N., Schweitzer, P., Van Leeuwen, E.J., Mehlhorn, K. and Borgwardt, K.M., 2011. Weisfeiler-lehman graph kernels. *Journal of Machine Learning Research*, 12(9).
- [5] Chen, Z., Villar, S., Chen, L. and Bruna, J., 2019. On the equivalence between graph isomorphism testing and function approximation with gns. *Advances in neural information processing systems*, 32.
- [6] Brül Gabrielsson, R., 2020. Universal Function Approximation on Graphs. *Advances in Neural Information Processing Systems*, 33, pp.19762-19772.
- [7] Loukas, A., 2019. What graph neural networks cannot learn: depth vs width. *arXiv preprint arXiv:1907.03199*.
- [8] Loukas, A., 2020. How hard is to distinguish graphs with graph neural networks?. *Advances in neural information processing systems*, 33, pp.3465-3476.
- [9] Chen, Z., Chen, L., Villar, S. and Bruna, J., 2020. Can graph neural networks count substructures?. *Advances in neural information processing systems*, 33, pp.10383-10395.
- [10] McAllester, D., 2003. Simplified PAC-Bayesian margin bounds. In *Learning theory and Kernel machines* (pp. 203-215). Springer, Berlin, Heidelberg.

References

- [11] Neyshabur, B., Bhojanapalli, S. and Srebro, N., 2017. A pac-bayesian approach to spectrally-normalized margin bounds for neural networks. arXiv preprint arXiv:1707.09564.
- [12] Kipf, T.N. and Welling, M., 2016. Semi-supervised classification with graph convolutional networks. arXiv preprint arXiv:1609.02907.
- [13] Dai, H., Dai, B. and Song, L., 2016, June. Discriminative embeddings of latent variable models for structured data. In International conference on machine learning (pp. 2702-2711). PMLR.
- [14] Liao, R., Urtasun, R. and Zemel, R., 2020. A pac-bayesian approach to generalization bounds for graph neural networks. arXiv preprint arXiv:2012.07690.
- [15] Garg, V., Jegelka, S. and Jaakkola, T., 2020, November. Generalization and representational limits of graph neural networks. In International Conference on Machine Learning (pp. 3419-3430). PMLR.
- [16] Scarselli, F., Tsoi, A.C. and Hagenbuchner, M., 2018. The vapnik–chervonenkis dimension of graph and recursive neural networks. *Neural Networks*, 108, pp.248-259.
- [17] Maron, H., Ben-Hamu, H., Shamir, N. and Lipman, Y., 2018. Invariant and equivariant graph networks. arXiv preprint arXiv:1812.09902.
- [18] Maron, H., Ben-Hamu, H., Serviansky, H. and Lipman, Y., 2019. Provably powerful graph networks. *Advances in neural information processing systems*, 32.
- [19] Maron, H., Ben-Hamu, H. and Lipman, Y., 2019. Open problems: Approximation power of invariant graph networks. In *NeurIPS 2019 Graph Representation Learning Workshop*, 2019a. URL <https://grlearning.github.io/papers/31.pdf> (Vol. 2, No. 8, p. 31).

Questions?