

CS480/680: Introduction to Machine Learning

Lec 10: Graph Neural Network

Yaoliang Yu

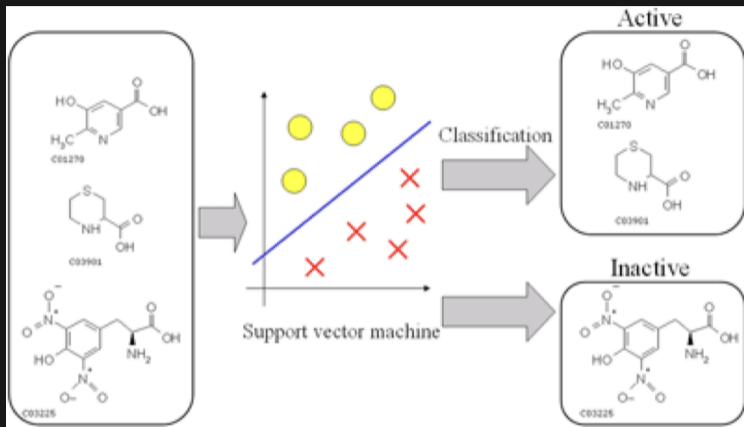


UNIVERSITY OF
WATERLOO

FACULTY OF MATHEMATICS
**DAVID R. CHERITON SCHOOL
OF COMPUTER SCIENCE**

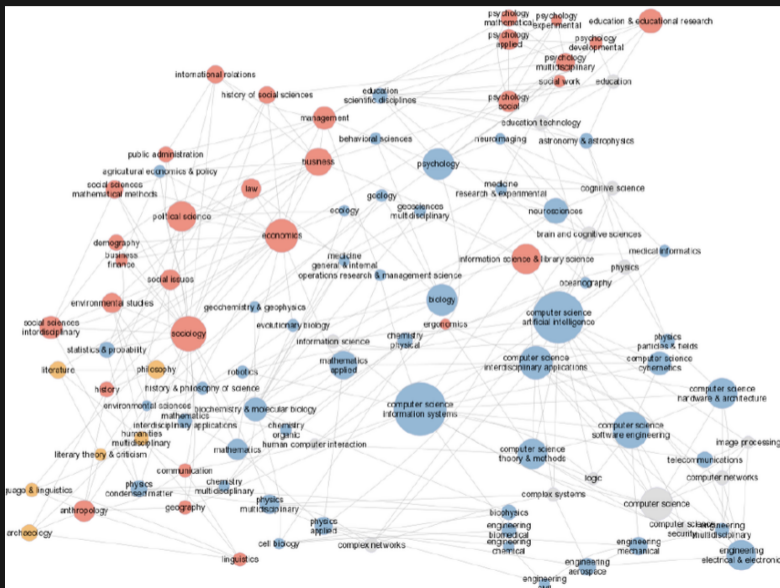
Feb 11, 2025

Chemical Compound



- Nodes are not necessarily in correspondence
- Output: e.g., function mapping unseen compound to level of activity against cancer cells

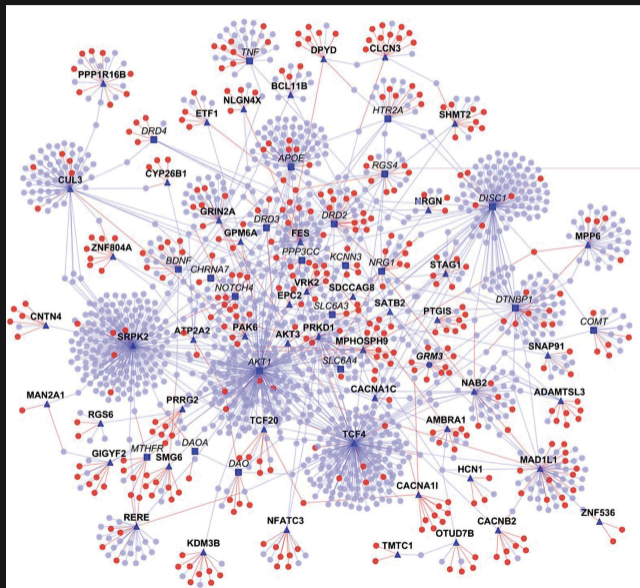
Collaboration Network



Social Network



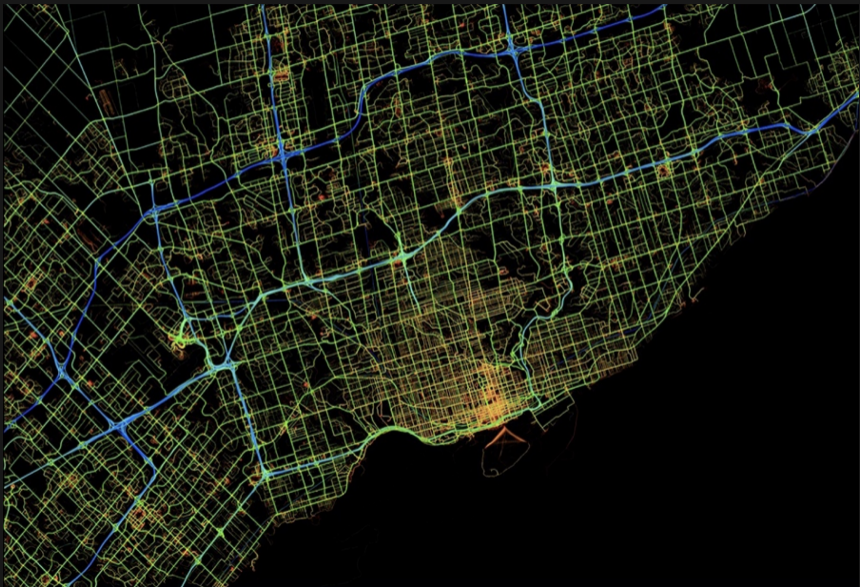
Biological Network



Communication Network

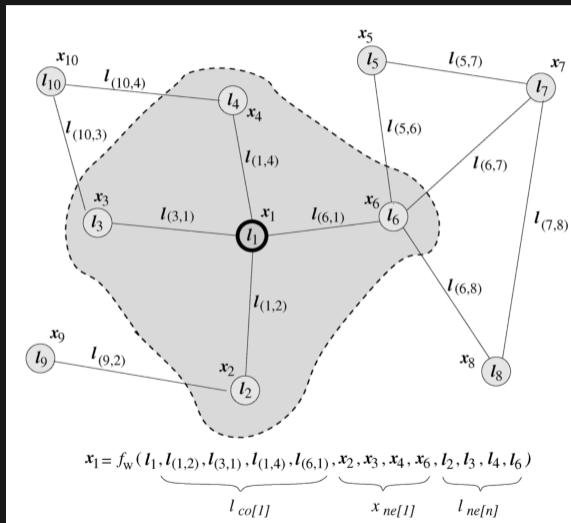


Traffic Network



Graph Neural Networks (GNN)

- Input: $\mathcal{G} = (\mathcal{V}, \mathcal{E}, l)$
- State (embedding):
$$h_u = f_{\mathbf{w}}(h_{ne[u]}, l_{ne[u]}, l_{co[u]})$$
- Output: $o_u = g_{\mathbf{w}}(h_u, l_u)$
- f, g : neural nets
- Parameters: \mathbf{w} , **shared among nodes**



State Update

$$\left. \begin{aligned} h_1 &= f_{\mathbf{w}}(h_{ne[1]}, l_{ne[1]}, l_{co[1]}) \\ &\vdots \\ h_n &= f_{\mathbf{w}}(h_{ne[n]}, l_{ne[n]}, l_{co[n]}) \end{aligned} \right\} \implies \mathbf{h} = F_{\mathbf{w}}(\mathbf{h}, \mathfrak{l})$$

- If $F_{\mathbf{w}}$ is a **contraction**, then exists a unique state \mathbf{h} :

$$\|F_{\mathbf{w}}(\mathbf{h}, \mathfrak{l}) - F_{\mathbf{w}}(\mathbf{z}, \mathfrak{l})\| \leq \gamma \|\mathbf{h} - \mathbf{z}\|, \text{ for some } \gamma \in [0, 1)$$

- $\mathbf{h}^{t+1} \leftarrow F_{\mathbf{w}}(\mathbf{h}^t, \mathfrak{l})$ converges linearly to the fixed point \mathbf{h}
- Upon convergence, output $\mathbf{o} = G_{\mathbf{w}}(\mathbf{h}, \mathfrak{l}) =: \hat{\mathbf{y}}(\mathfrak{l}; \mathbf{w})$

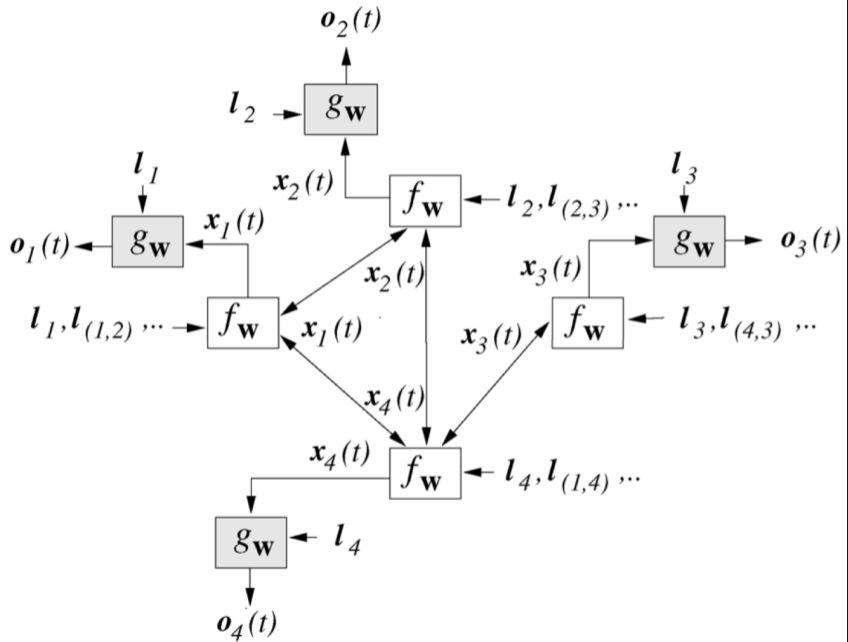
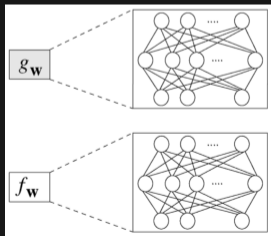
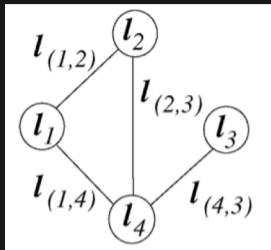
Example

$$h_v = \frac{1}{|ne[v]|} \sum_{u \in ne[v]} \varphi_{\mathbf{w}}(h_u, l_v, l_{(v,u)}, l_u)$$

- State h_u , node feature l_u and the edge feature $l_{(v,u)}$ can all be vectors
- Graph structure is used in the sum: only neighbors contribute
- PageRank:

$$h_v = \sum_{u \in ne[v]} a_{vu} h_u, \quad \text{e.g. } a_{vu} := \frac{1}{|ne[v]|} \mathbb{1}[v \in ne[u]]$$

- In GNN, the aggregation function $\varphi_{\mathbf{w}}$ is learnable through \mathbf{w}



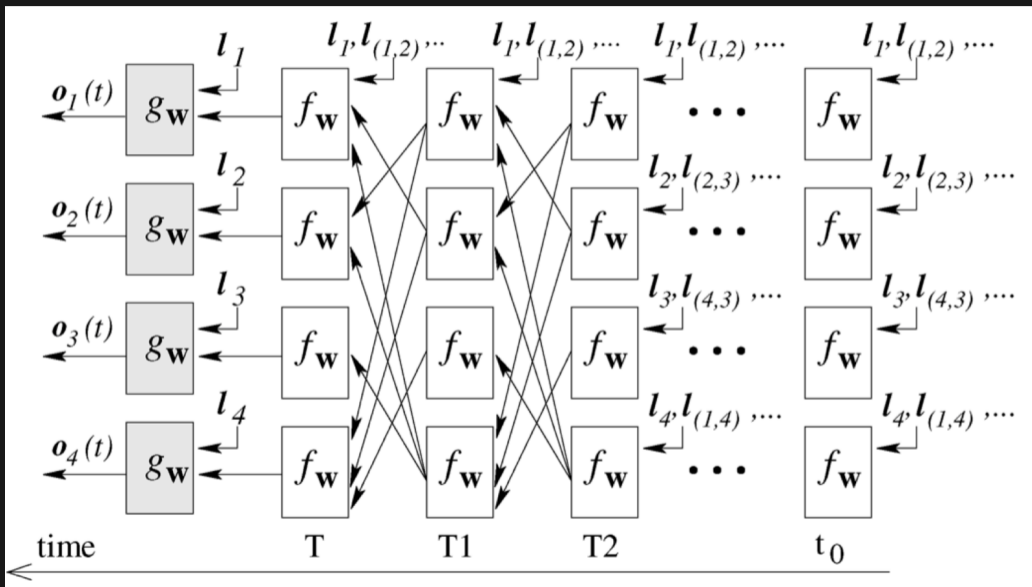
Learning GNN

- Given a supervised set of graphs and labels $(\mathcal{G}_i, \mathbf{y}_i)$, $i = 1, \dots, n$
- Learn a predictor $\hat{\mathbf{y}}$ that maps a new test graph \mathcal{G} to its label: $\hat{\mathbf{y}}(\mathcal{G}) \approx \mathbf{y}$
 - labels could be at the node, edge or graph level
 - do not confuse the label \mathbf{y} with the feature \mathbf{l}
- Choose a loss function L to solve:

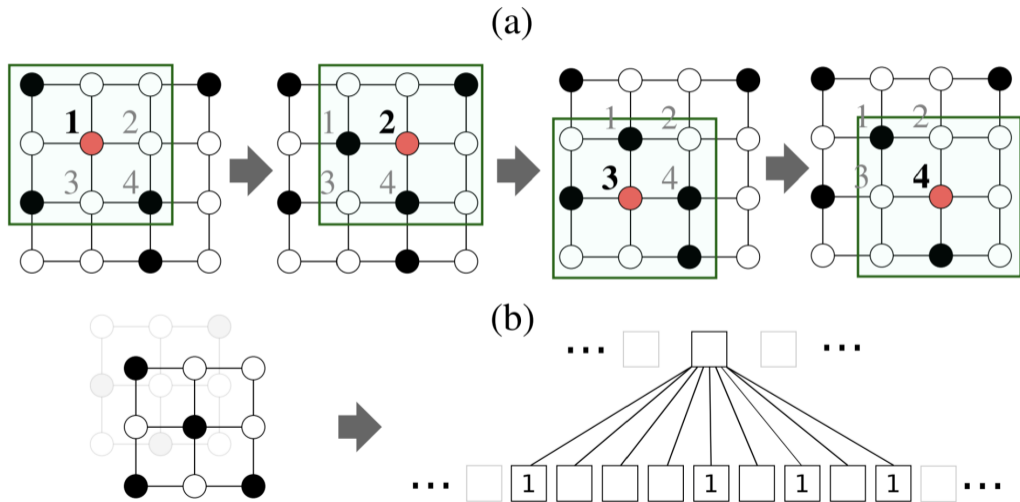
$$\min_{\mathbf{w}} L(\hat{\mathbf{y}}(\mathbf{l}; \mathbf{w}), \mathbf{y})$$

- unroll k steps: $\hat{\mathbf{y}}(\mathbf{l}; \mathbf{w}) \approx G_{\mathbf{w}}(F_{\mathbf{w}}^{[k]}(\mathbf{h}, \mathbf{l}), \mathbf{l})$
- or apply implicit function theorem to differentiate \mathbf{w}

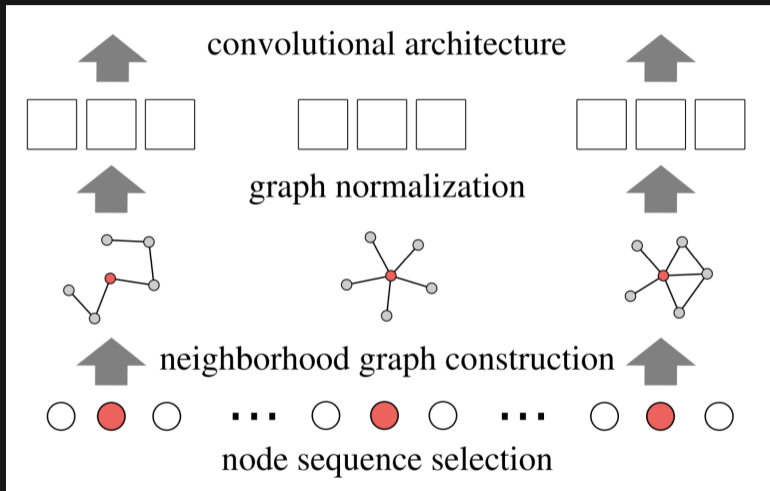
Training by Backpropagation Through Time (BPTT)



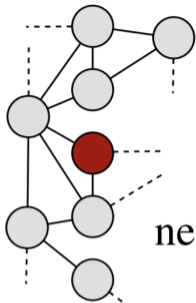
CNN Recalled



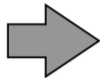
Spatial Convolution



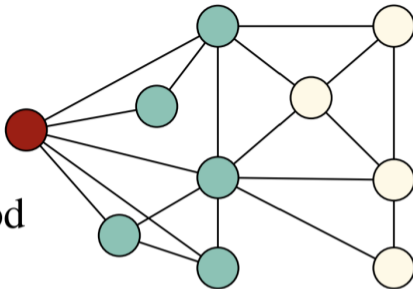
M. Niepert, M. Ahmed, and K. Kutzkov. "Learning Convolutional Neural Networks for Graphs". In: *Proceedings of The 33rd International Conference on Machine Learning*. 2016, pp. 2014–2023.



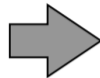
select



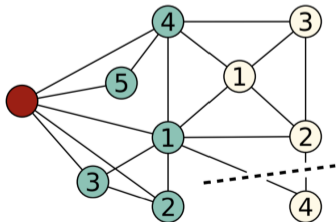
neighborhood



normalize



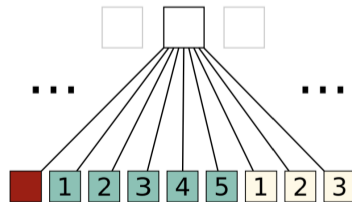
subgraph



receptive field



reads vertex and edge attributes = channels



Fourier Transform

$$f * g = \mathcal{F}^{-1} (\mathcal{F}[f] \cdot \mathcal{F}[g])$$

- Convolution in time domain = multiplication in frequency domain
- Invertible, in fact, orthogonal transform
- Can we do something similar for graphs?

Graph Laplacian

$$A_{uv} = \mathbb{I}[(u, v) \in \mathcal{E}], \quad D_{uv} = \sum_n A_{un} \cdot \mathbb{I}[u = v]$$

- Laplacian $L = D - A$
 - $L\mathbf{1} = 0 \cdot \mathbf{1}$: # connected components of \mathcal{G} = # multiplicity of $\lambda = 0$
 - symmetric and PSD for undirected graph: $\langle \mathbf{x}, L\mathbf{x} \rangle = \frac{1}{2} \sum_{u,v} A_{uv} (x_u - x_v)^2$
- Normalized Laplacian $\bar{L} = I - D^{-1/2}AD^{-1/2} = D^{-1/2}LD^{-1/2}$

Example

$$A = \begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix}, \quad D = \begin{bmatrix} 2 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 2 \end{bmatrix} = \text{diag}([2, 2, 2])$$

$$L = D - A = \begin{bmatrix} 2 & -1 & -1 \\ -1 & 2 & -1 \\ -1 & -1 & 2 \end{bmatrix}, \quad \bar{L} = D^{-1/2} L D^{-1/2} = \begin{bmatrix} 1 & -\frac{1}{2} & -\frac{1}{2} \\ -\frac{1}{2} & 1 & -\frac{1}{2} \\ -\frac{1}{2} & -\frac{1}{2} & 1 \end{bmatrix}$$

$$\mathring{A} = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}, \quad \mathring{D} = \begin{bmatrix} 3 & 0 & 0 \\ 0 & 3 & 0 \\ 0 & 0 & 3 \end{bmatrix}$$

Spectral Convolution

- Given two graph signals $\mathbf{x}, \mathbf{g} \in \mathbb{R}^{|\mathcal{V}|}$:

$$\mathbf{x} * \mathbf{g} := U[(U^\top \mathbf{x}) \odot (U^\top \mathbf{g})], \quad \text{where } L = U \Lambda U^\top$$

- One layer of spectral convolution, with filter weights W_r^k to be learned:

$$\mathbf{x}_r^{k+1} := \sigma \left(\underbrace{U[W_r^k \odot (U^\top X^k)] \mathbf{1}}_{\text{aggregation along depth}} \right), \quad r = 1, \dots, d_{k+1}, \quad X^k = [\mathbf{x}_1^k, \dots, \mathbf{x}_{d_k}^k]$$

conv by r -th filter

- Can stack to go deep

Chebyshev Net

- Spectral conv requires eigen-decomposition and is not localized
- Rewrite the convolution:

$$\begin{aligned}\mathbf{x} * \mathbf{g} &:= U[(U^\top \mathbf{x}) \odot (U^\top \mathbf{g})] = U[\text{diag}(f(\lambda; \mathbf{w}))(U^\top \mathbf{x})] \\ &= [U \text{diag}(f(\lambda; \mathbf{w}))U^\top] \mathbf{x} \\ &:= f(L; \mathbf{w}) \mathbf{x}\end{aligned}$$

- Choosing f to be a polynomial dispenses eigen-decomposition
- Resulting conv is localized: degree k polynomial only requires k -hop neighbors

Graph Convolutional Net (GCN)

- A layer of GCN is defined simply as:

$$X^{k+1} = \sigma(\mathring{D}^{-1/2} \mathring{A} \mathring{D}^{-1/2} X^k W^k), \quad X^k = [\mathbf{x}_1^k, \dots, \mathbf{x}_s^k] \in \mathbb{R}^{|\mathcal{V}| \times s}, \quad W^k \in \mathbb{R}^{s \times t}$$

- $\mathring{A} = A + I$ (i.e. adding self-cycle)
 - \mathring{D} is the usual diagonal degree matrix of \mathring{A}
 - s and t are the number of input and output channels, resp.
- One layer of GCN only aggregates info from 1-hop neighbors
 - Can stack to get deep and aggregate info from k -hop neighbors

Connections

- Rewriting GCN in vector form and identify $X_v = \mathbf{l}_v$:

$$\mathbf{l}_v^{k+1} = \sigma \left(\left[\frac{1}{d_v + 1} \mathbf{l}_v^k + \sum_{u \in \mathcal{N}_v} \frac{a_{vu}}{\sqrt{(d_v + 1)(d_u + 1)}} \mathbf{l}_u^k \right] W^k \right)$$

- This is GNN!
- It also resembles the Weisfeiler-Lehman (WL) algorithm!

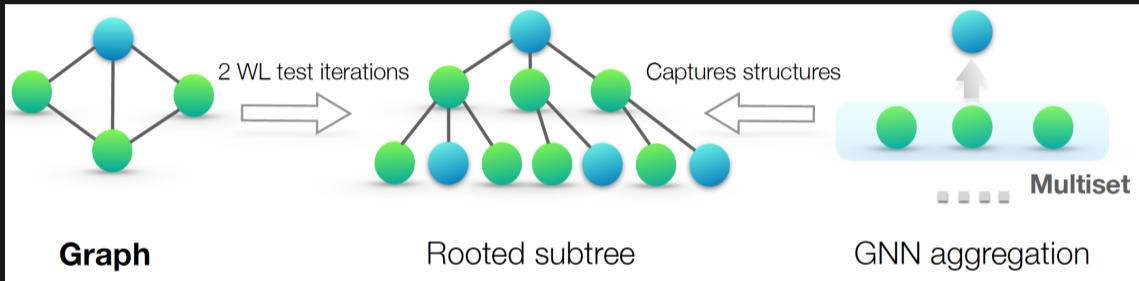
WL and Graph Isomorphism

Algorithm 1: Weisfeiler-Lehman iterative color refinement

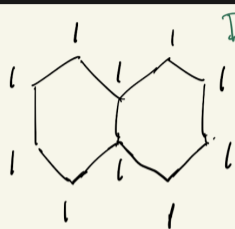
Input: Graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \iota^0)$

Output: $\iota^{|\mathcal{V}|-1}$

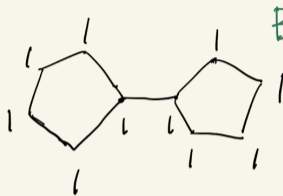
```
1 for  $t = 0, 1, \dots, |\mathcal{V}| - 1$  do
2    $\iota^{t+1} \leftarrow \text{hash}([\iota_v^t, \iota_{u \in \mathcal{N}_v}^t] : v \in \mathcal{V})$  //  $[\cdot]$  is a multiset, allowing repetitions
3 Function  $\text{hash}([\iota_v, \iota_{u \in \mathcal{N}_v}] : v \in \mathcal{V})$ :
4   for  $v \in \mathcal{V}$  do
5      $\text{sort}(\iota_{u \in \mathcal{N}_v})$  // sort the neighbors
6     prefix  $\iota_v$  to sorted list  $[\iota_v, \iota_{u \in \mathcal{N}_v}]$  //  $\iota_v$  does not participate in sorting!
7      $\iota_v^+ \leftarrow f([\iota_v, \iota_{u \in \mathcal{N}_v}])$  //  $f : L^* \rightarrow L$  lexicographic strictly increasing
```



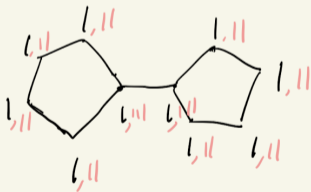
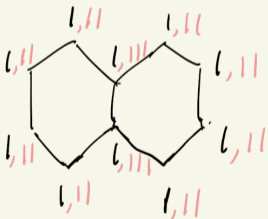
Example



Decalin

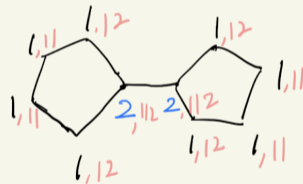
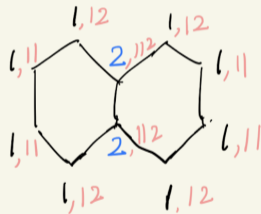
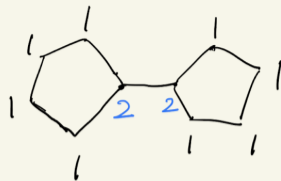
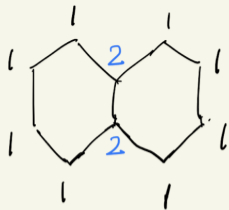


Bicyclopentyl



$$f: (1, II) \rightarrow 1$$

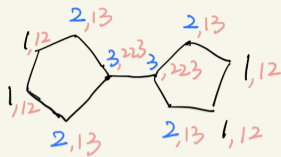
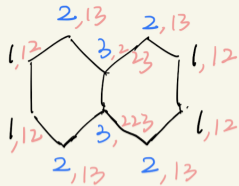
$$(1, III) \rightarrow 2$$



$$f: (1,11) \rightarrow 1$$

$$(1,12) \rightarrow 2$$

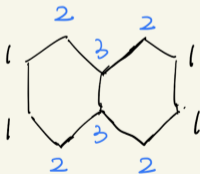
$$(2,112) \rightarrow 3$$



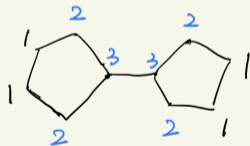
$$f: (1, 12) \rightarrow 1$$

$$(2, 13) \rightarrow 2$$

$$(3, 223) \rightarrow 3$$



\neq



Converged!

WL cannot distinguish the two graphs.

From 1 to 2

- GCN: $l_v \leftarrow \sigma \left(\underbrace{\left[\frac{1}{d_v + 1} l_v + \sum_{u \in \mathcal{N}_v} \frac{a_{vu}}{\sqrt{(d_v + 1)(d_u + 1)}} l_u \right]}_{\text{averaged input}} W \right)$
linear layer

- GraphSAGE: $l_v \leftarrow \sigma \left(\underbrace{\left[l_v \vee \max_{u \in \mathcal{N}_v} \{l_u\} \right]}_{\text{max-pooled input}} W \right)$
linear layer

Graph Isomorphism Network (GIN)

$$\mathbf{l}_v \leftarrow \text{MLP} \left(\underbrace{(1 + \epsilon)\mathbf{l}_v + \sum_{u \in \mathcal{N}(v)} \mathbf{l}_u}_{\text{summed input}} \right)$$

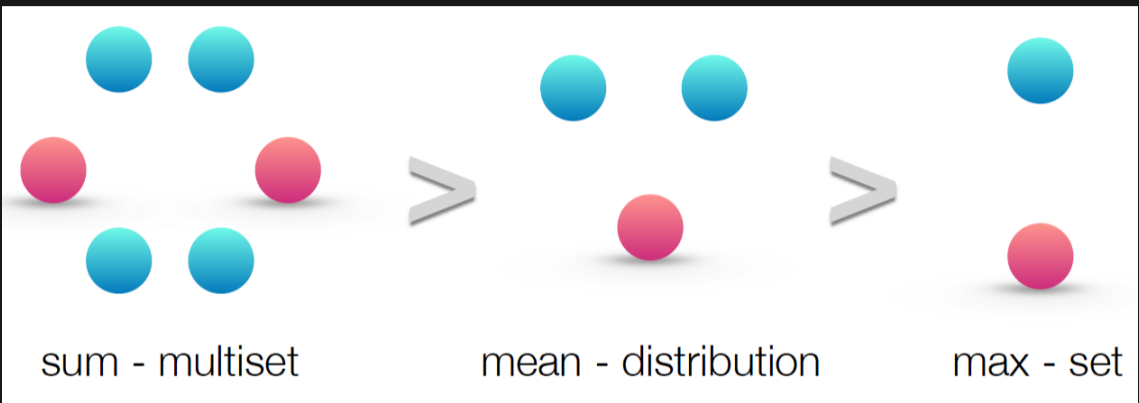
Theorem: Representation power of GNNs

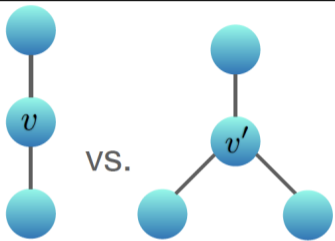
GNNs with aggregation function

$$\mathbf{l}_v \leftarrow \sigma \left(\mathbf{l}_v, \varphi(\mathbf{l}_u : u \in \mathcal{N}(v)) \right).$$

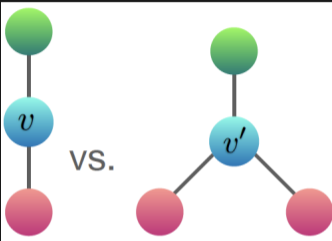
are no more discriminative than WL, with equality if φ and σ are injective.

- Compared to WL, GNNs also tend to map similar nodes into similar embeddings

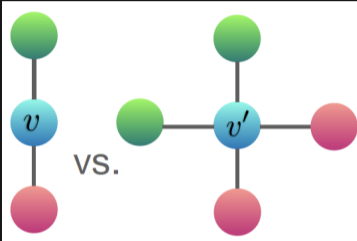




(a) Mean and Max both fail



(b) Max fails



(c) Mean and Max both fail

