

This notes not for newbie
 (not just know how but know why)

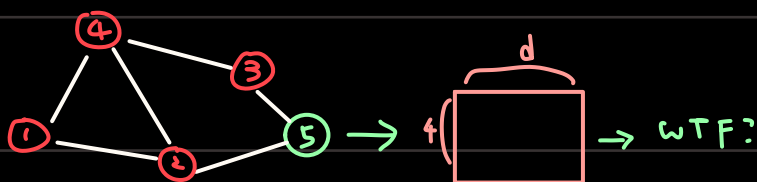
Assump you know: DL for CV, optimization, LSTM, RNN, Attention

Deep Learning For Graphs

Weak of node embedding:

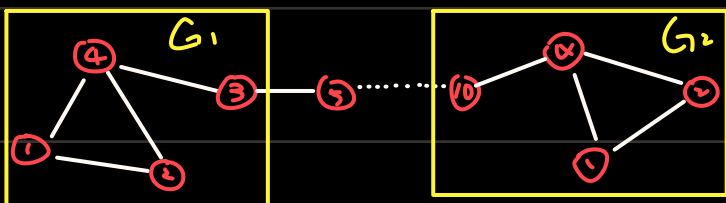
①

lookup-table = limitation



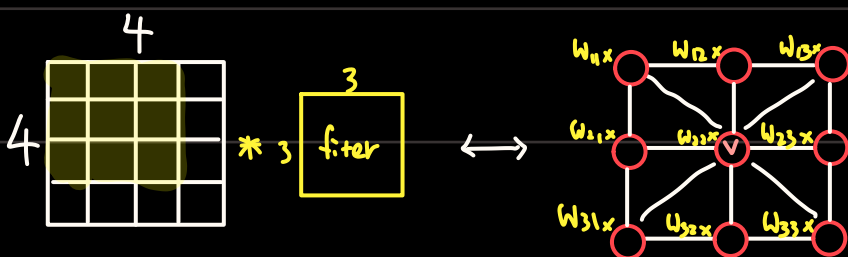
②

Structure similarity: in space they are not close in big graph

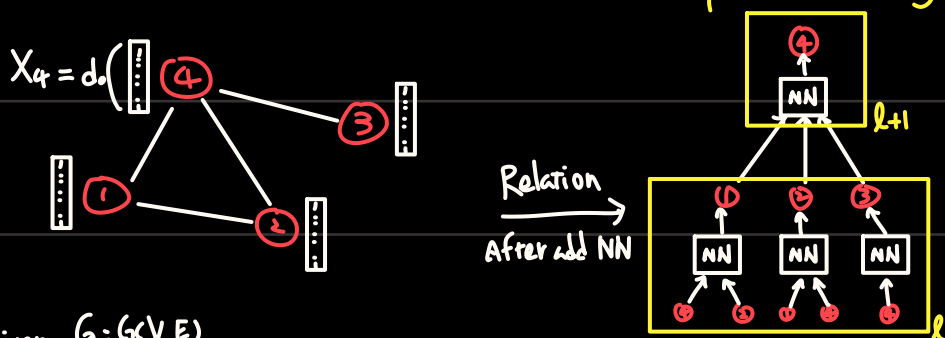


Def Aggregate from Neighbors (Basic Graph NN: NN4Graph)

Image is a tidy graph \rightarrow CNN use pixel itself and its neighbors to transform a new pixel (feature mapping)



How about we use a shared weight Matrix to all the graph?
 (since in CNN use 1 filter to do all convolution operation in a single layer)



Given $G=(V,E)$

(Forward):

For each node $v \in V$ with node feature $h_v^{(0)} = x_v = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$ do

$$\{ h_v^{(l+1)} = \sigma \left(w^{(l)} \sum_{u \in N(v)} \frac{h_u^{(l)}}{|N(v)|} + b^{(l)} \cdot h_v^{(l)} \right), \forall l \in \{0, \dots, L-1\} \}$$

$\rightarrow z_v = h_v^{(L)} \in \mathbb{R}^d$ = final node embedding

(for loop will take too much time complexity to do single forward)

Matrix Formulation:

Let $H^{(l)} = [h_1^{(l)} \dots h_{N(v)}^{(l)}]^T \rightarrow \sum_{u \in N(v)} h_u^{(l)} = A_v \cdot H^{(l)}$, where A_v is no. v row in adjacency matrix

Let D be a diagonal matrix where $D_{v,v} = \text{Deg}(v) = |N(v)| \rightarrow D_{v,v}^{-1} = \frac{1}{|N(v)|}$

$\rightarrow \sum_{u \in N(v)} \frac{h_u^{(l)}}{|N(v)|} \forall v \in V \Rightarrow D^{-1} \cdot A \cdot H^{(l)}$ where $D^{N \times N}$, $A^{N \times N}$, $H^{N \times d}$, $d \in \{d_0, \dots, d\}$

$$\rightarrow \begin{bmatrix} \frac{1}{|N(v_1)|} & & 0 \\ & \ddots & \\ 0 & & \frac{1}{|N(v_N)|} \end{bmatrix} \left(\begin{bmatrix} \text{Adj} \in \{0,1\}^{N \times N} \end{bmatrix} \begin{bmatrix} h_1^{(l)} \\ \vdots \\ h_N^{(l)} \end{bmatrix} \right) \text{Neighbors}$$

$\rightarrow h_v^{(l+1)} = \sigma \left(\omega^{(l)} \sum_{u \in N(v)} \frac{h_u^{(l)}}{|N(v)|} + b^{(l)} \cdot h_v^{(l)} \right) \Rightarrow H^{(l+1)} = \sigma \left(D^{-1} A \cdot H^{(l)} \cdot \omega^{(l)T} + H^{(l)} \cdot b^{(l)T} \right)$

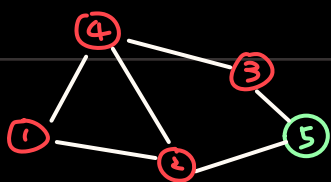
$\rightarrow \tilde{Z} = \begin{pmatrix} z_1 \\ \vdots \\ z_{N(v)} \end{pmatrix} = H^{(L)} \in \mathbb{R}^{d \times N(v)}$

train \rightarrow Unsupervised: $\min_{\omega, b} \left\| \text{Adj} - \frac{\tilde{Z}^T \cdot \tilde{Z}}{\|\tilde{Z}\| \|\tilde{Z}\|} \right\|_F^2$
 low bias norm

Supervised: Node classification (Use linear layer + sigmoid then do backprop to update weights)

Previous limitation discussion:

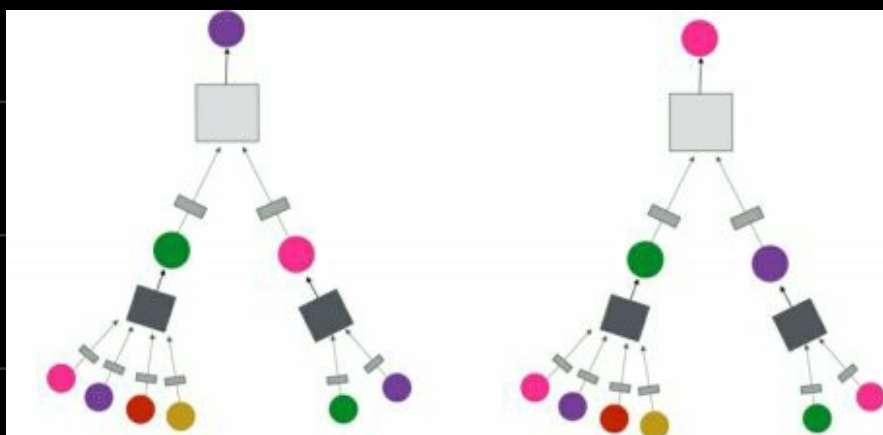
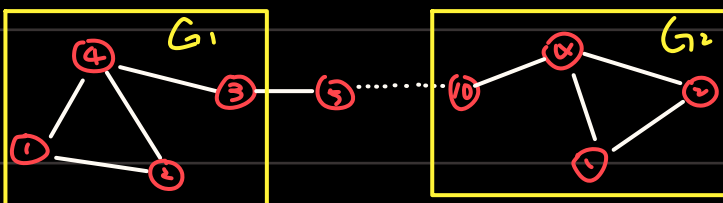
①



$$X_5 = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{pmatrix} \Rightarrow h_5^{(0)} = h_5^{(0)} \rightarrow h_v^{(l+1)} = \sigma \left(\omega^{(l)} \sum_{u \in N(v)} \frac{h_u^{(l)}}{|N(v)|} + b^{(l)} \cdot h_v^{(l)} \right) \dots \rightarrow h_5^{(L)} \text{ (solved)}$$

$\omega^{(L)} \begin{bmatrix} h_1^{(L)} \\ h_2^{(L)} \\ h_3^{(L)} \end{bmatrix} \cdot \frac{1}{2} + b^{(L)} \cdot h_5^{(L)}$

②

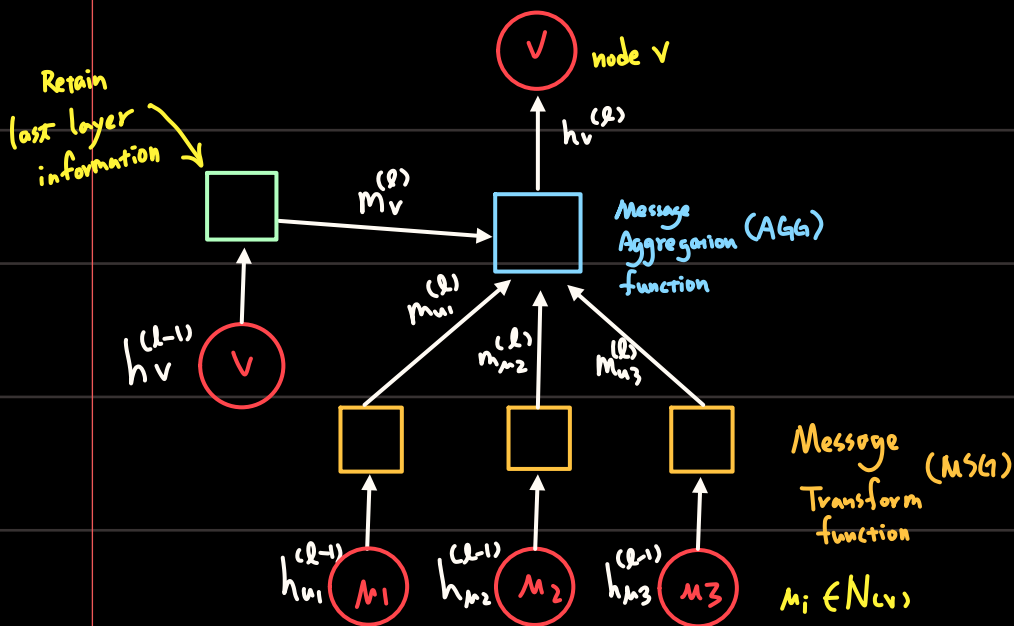


(solved)

Need to define

- ① A GNN layer = Message + Aggregation (GCN, GraphSage, GAT... differ these 2 components)
- ② How layer connect?
- ③ Graph augmentation, Structure manipulation
- ④ Task? How to train

(Def) Define A Single Layer = MSG + AGG (General Single Forward Formulation)



① Message Transformation

$$m_v^{(l)} = \text{MSG}^{(l)}(h_{v_i}^{(l-1)})$$

Many choices e.g. $w_{vi}^{(l)} \cdot h_{v_i}^{(l-1)}$ + Activation (or not)

② Message Aggregation

$$h_v^{(l)} = \text{AGG}(\{m_{v_i}^{(l)}, m_{v_i}^{(l-1)}\})$$

e.g. Sum, Mean, Max + Activation (or not)

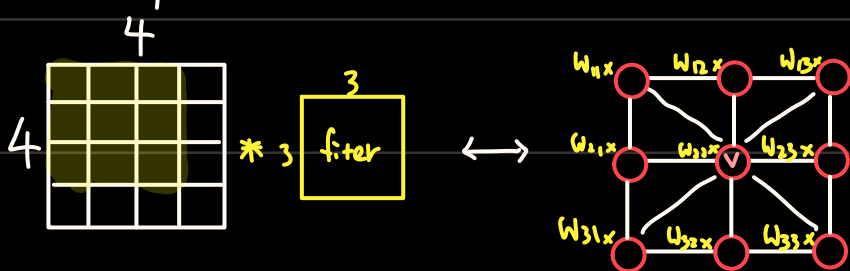
Issue (similar to ResNet, DenseNet): Add $h_v^{(l-1)}$ into message passing operation

$$\text{① } m_v^{(l)} = \beta \cdot h_v^{(l-1)}$$

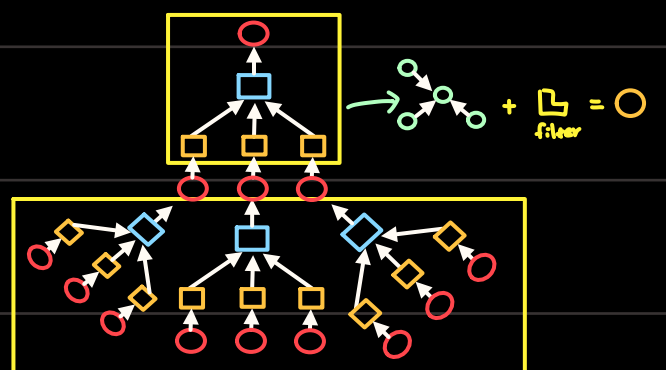
$$\text{② } h_v^{(l)} = \text{Concat}(\text{AGG}(\{m_{v_i}^{(l)}, m_{v_i}^{(l-1)}\}), m_v^{(l)})$$

(or Add)

(Def) Graph Convolutional Network (Spatial-Based Convolutional)



(Not like classic CNN)



Origin Basic Operator: $h_v^{(l)} = \sigma(w^{(l)} \sum_{u \in N(v)} \frac{h_u^{(l-1)}}{|N(v)|} + \beta \cdot h_v^{(l-1)})$

(NN+Graph)

$$\text{Message} \xrightarrow{+} \text{Aggregation} \rightarrow h_v^{(l)} = \sigma \left(\sum_{u \in N(v)} w^{(l)} \frac{h_u^{(l-1)}}{|N(v)|} + \beta \cdot h_v^{(l-1)} \right)$$

$$\begin{bmatrix} h_1^{(l+1)} \\ \vdots \\ h_{N(v)}^{(l+1)} \end{bmatrix} = \sigma \left[\begin{bmatrix} 1 & \dots & 0 \\ 0 & \dots & 1 \\ \vdots & \vdots & \vdots \end{bmatrix} \begin{bmatrix} w^{(l)} \\ \vdots \\ w^{(l)} \end{bmatrix} \begin{bmatrix} \text{Adj} \leftarrow \text{so, } r_i^{(l-1)} \end{bmatrix} + \begin{bmatrix} \beta \\ \vdots \\ \beta \end{bmatrix} \right] \begin{bmatrix} h_1^{(l)} \\ \vdots \\ h_{N(v)}^{(l)} \end{bmatrix}$$

$$\text{MSG} = \text{Linear} + \text{Normalize Neighbors} = \frac{w^{(l)} \cdot h_u^{(l-1)}}{|N(v)|}, \forall u \in N(v)$$

$$\text{AGG} = \text{Sum} + \text{Activation} = \sigma(\text{Sum}(\{m_u^{(l)}, m_u^{(l-1)}\}))$$

Def GraphSAGE (similar to DenseNet : $\text{Conv}(\text{concat}([a_1, a_2, \dots]))$)

$$h_v^{(l)} = \sigma(W^{(l)} \text{Concat}([h_v^{(l-1)}, \text{AGG}(\{h_m^{(l-1)}, m \in N(v)\})]))$$

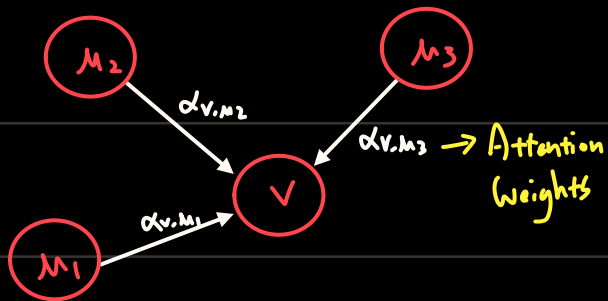
Stage I AGG choices :

- ① $\sum_{m \in N(v)} \frac{h_m^{(l-1)}}{\|N(v)\|} \text{MSG}$
- ② or Max $\text{Mean}(\{\text{MLP}(h_m^{(l-1)}), \forall m \in N(v)\})$ (pooling) MSG
- ③ Consider order : $\text{LSTM}([h_m^{(l-1)}, \forall m \in N(v)])$ No MSG

$$\text{Stage II AGG} : \sigma(W^{(l)} \cdot \text{Concat}(h_v^{(l-1)}, m_{m \in N(v)}^{(l)}))$$

Add l_2 -Normalization when pass : $h_v^{(l)} = \frac{h_v^{(l)}}{\|h_v^{(l)}\|_2} \quad \forall v, \forall l$

Def Graph Attention Networks



Which neighbors should node v need to pay attention?
(If ignore attention $\alpha_{v,m} = \alpha_{v,m_2} = \alpha_{v,m_3}$)

$$h_v^{(l)} = \sigma(\sum_{m \in N(v)} \alpha_{v,m} W^{(l)} h_m^{(l-1)}) \quad (\text{Concept like weighted sum})$$

Attention Coefficient e_{vu} :

e_{vu} = the importance of $m_m^{(l)}$ passing to v

$= a(W^{(l)} h_v^{(l-1)}, W^{(l)} h_m^{(l-1)})$, where a is a function which has many choices

$a(\cdot)$ choices : dot-product, concat + linear, ...

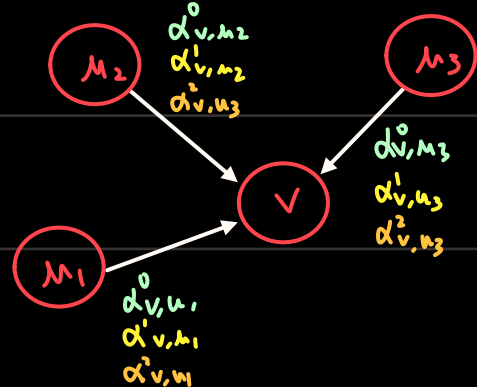
$$\text{Weight} : \alpha_{vu} = \frac{\exp(e_{vu})}{\sum_{k \in N(v)} \exp(e_{vk})} \quad (\text{softmax})$$

Multi-Head Attention (From Transformer Structure)

$$h_v^{(l)} = \text{AGG} \left(\begin{array}{l} h_v^{(l)[0]} = \sigma(\sum_{m \in N(v)} \alpha_{v,m}^0 W^{(l)} h_m^{(l-1)}) \\ h_v^{(l)[1]} = \sigma(\sum_{m \in N(v)} \alpha_{v,m}^1 W^{(l)} h_m^{(l-1)}) \\ h_v^{(l)[2]} = \sigma(\sum_{m \in N(v)} \alpha_{v,m}^2 W^{(l)} h_m^{(l-1)}) \end{array} \right)$$

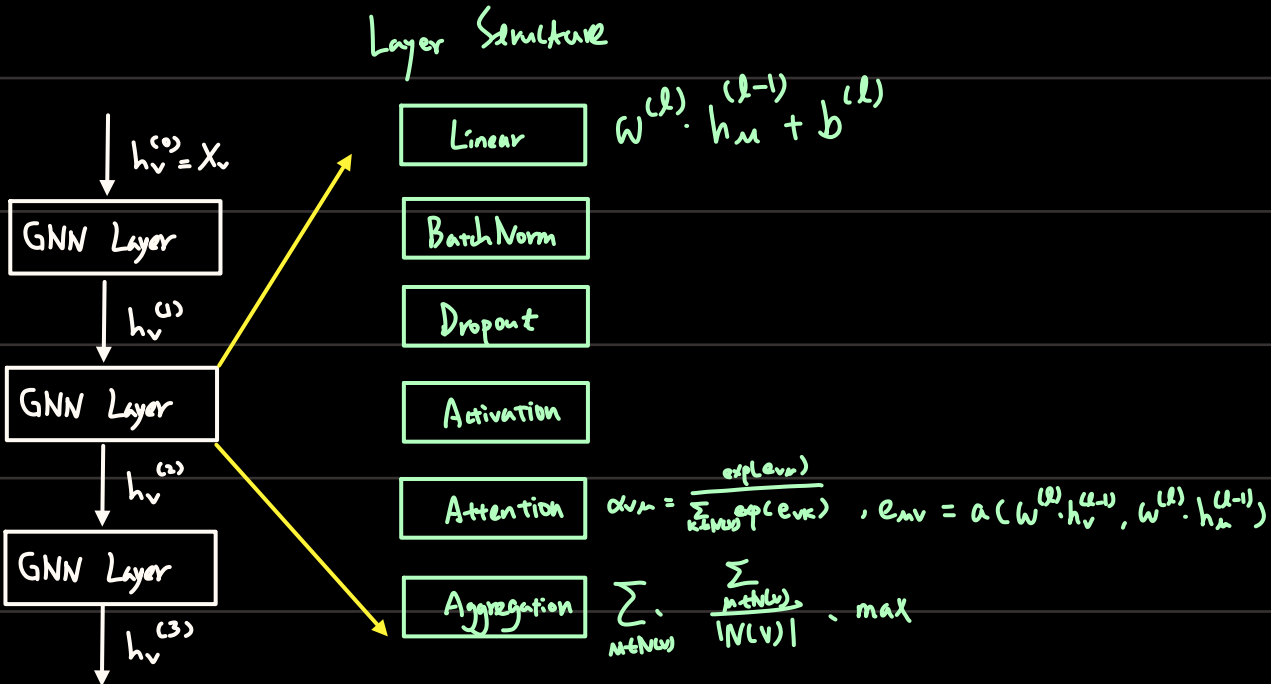
Sum-Concat

Example : 3 Attention-Head



Stacking GNN Layers

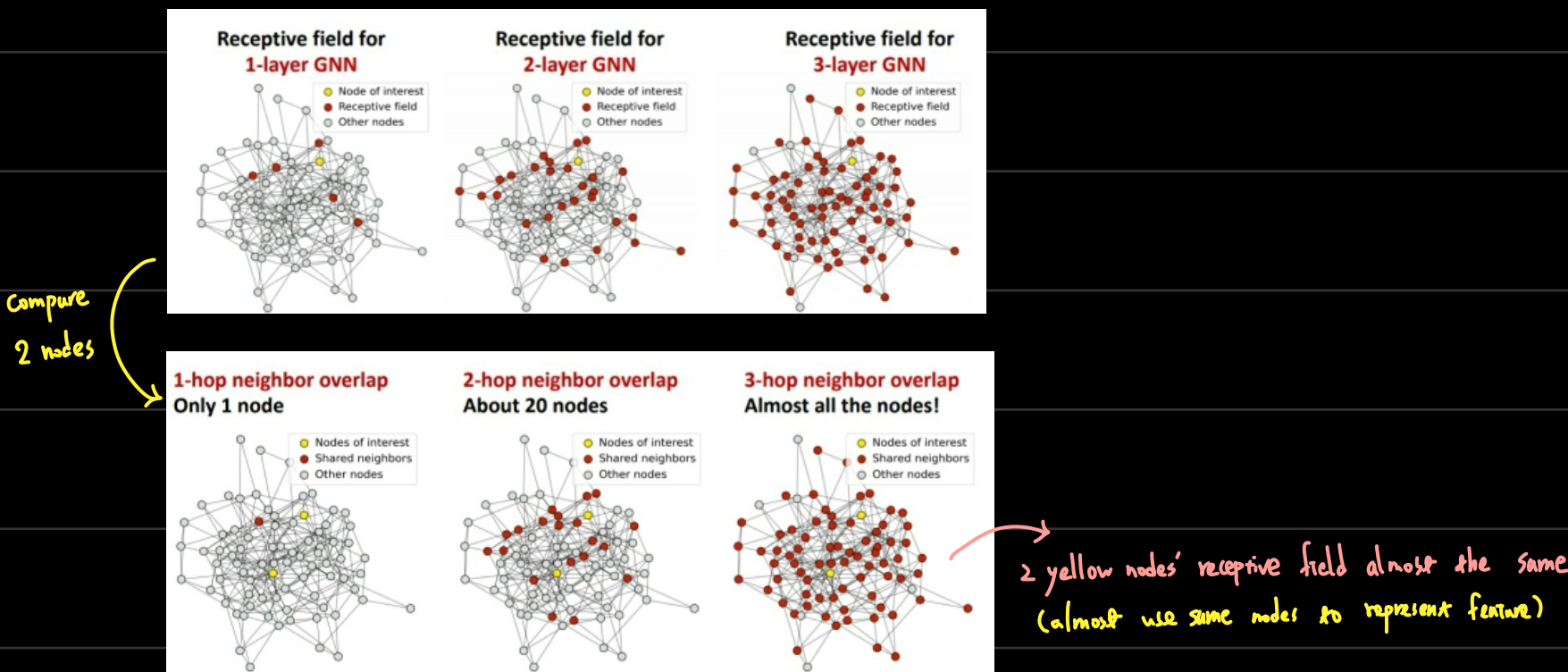
Construct Sequentially →



The Over-Smoothing Problem: Stack too many layers → all embeddings converge to the same value (or very similar)
 (similar to gradient vanishment) (all NN collect same information)

Receptive Field: the set of nodes that determine the embedding of a node of interest

(e.g. K-layer GNN: each node has a receptive field of K-hop neighborhood)



- OverCome:
- ① Analyze the necessary receptive field (e.g. compute the diameter of the graph)
 - ② Set num of GNN Layer L to be a bit more than receptive field

③ How to enhance expression with small layer GNN? → use MLP:

④ Add Layer that do not pass messages:

