

Graph Neural Network

Hyelin Choi

Department of Mathematics
Sungkyunkwan University

February 22, 2024

Table of Contents

- I. Convolutional Neighborhood Aggregation Method
 - I. Graph Neural Network
 - II. Graph Convolutional Network
 - III. GraphSAGE
 - IV. Graph Attention Network

Table of Contents

I. Convolutional Neighborhood Aggregation Method

I. Graph Neural Network

II. Graph Convolutional Network

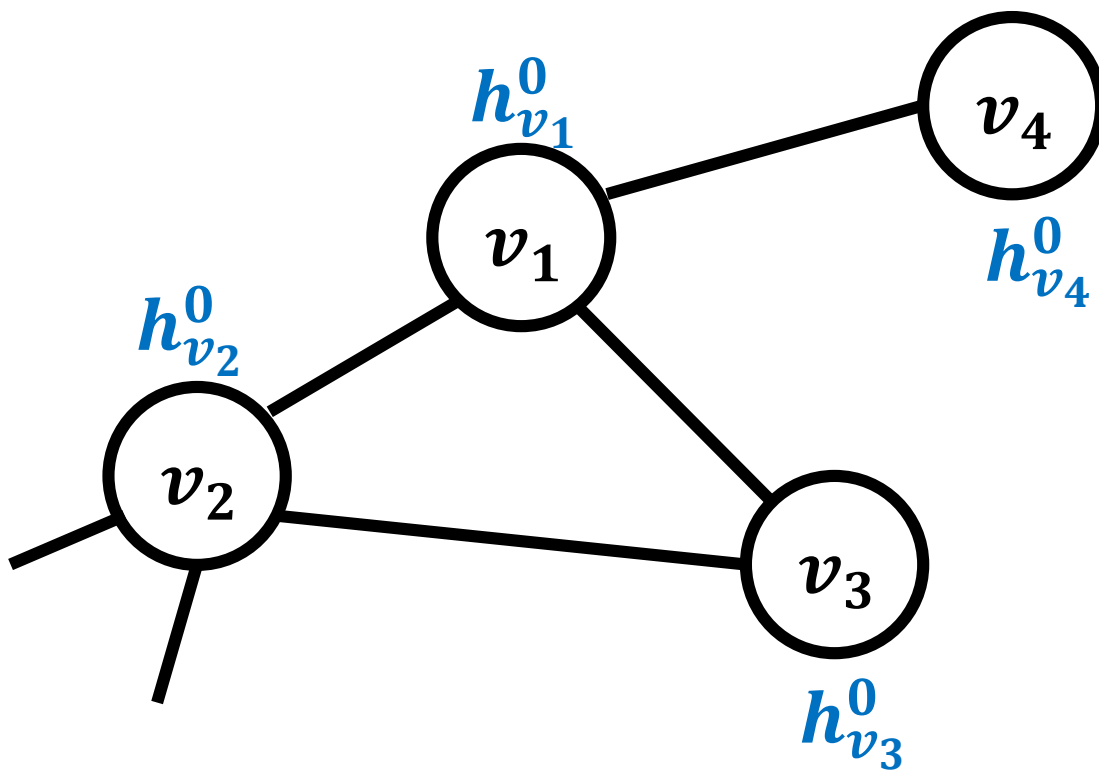
III. GraphSAGE

IV. Graph Attention Network

Convolutional Neighborhood Aggregation Method

Step 1

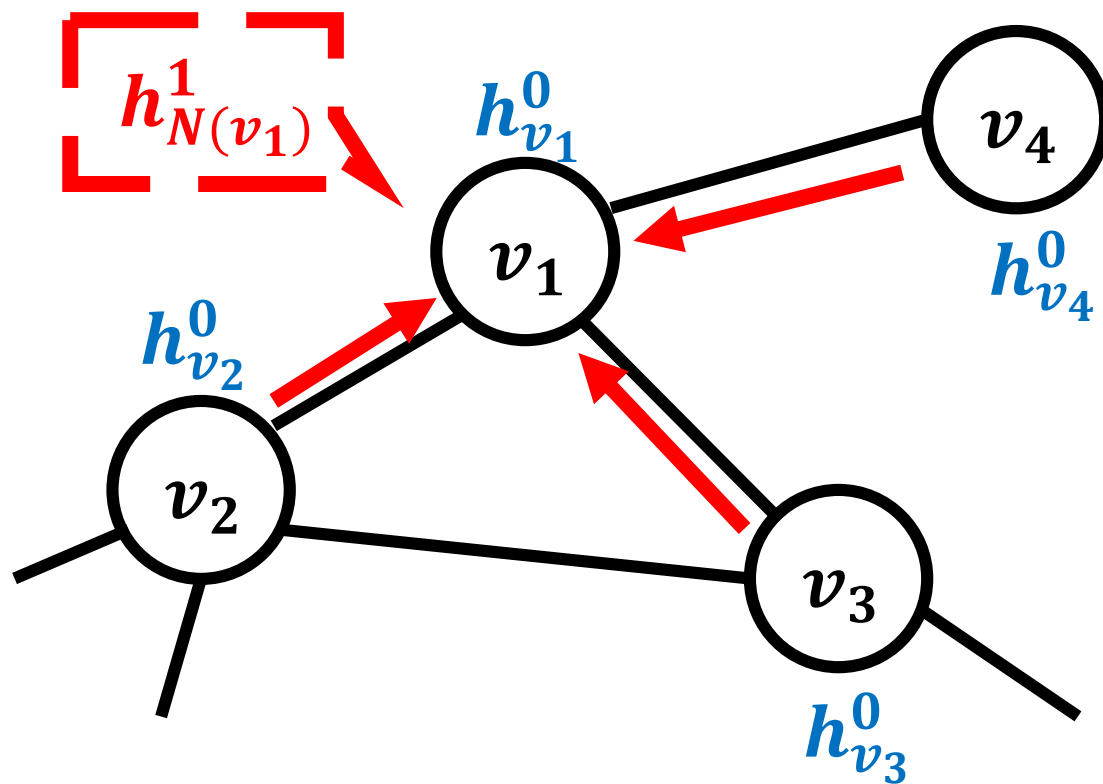
The node embeddings are initialized with node features.



Convolutional Neighborhood Aggregation Method

Step 2

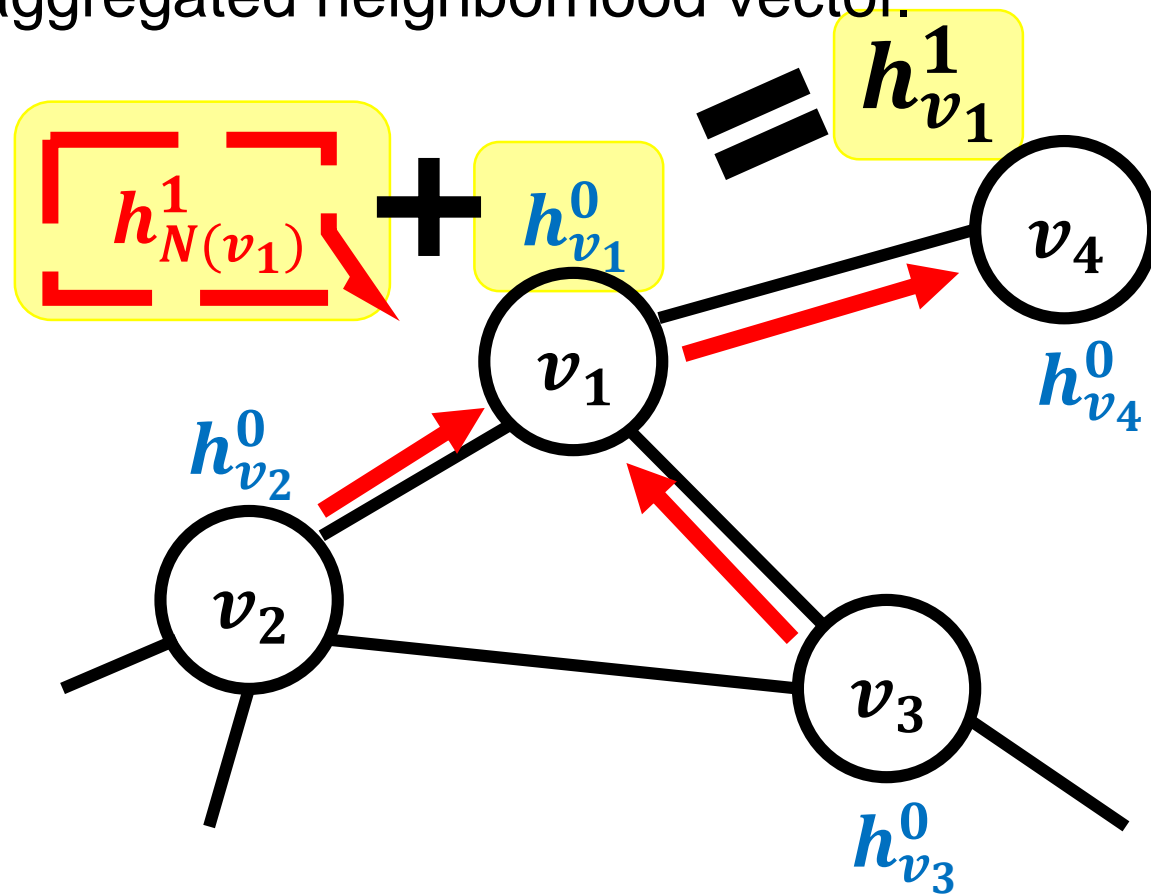
Nodes aggregate the embeddings of their neighbors, using an aggregation function.



Convolutional Neighborhood Aggregation Method

Step 3

Combine its previous embedding from the last iteration with its aggregated neighborhood vector.

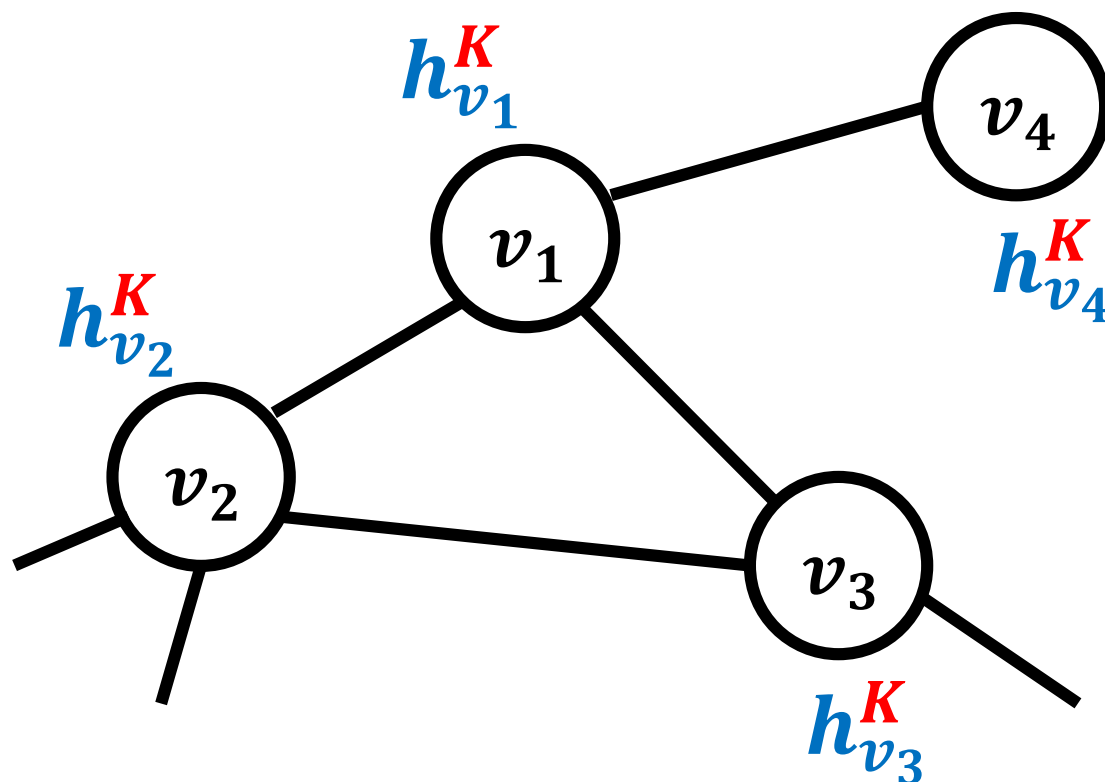


K iteration!

Convolutional Neighborhood Aggregation Method

Step 3

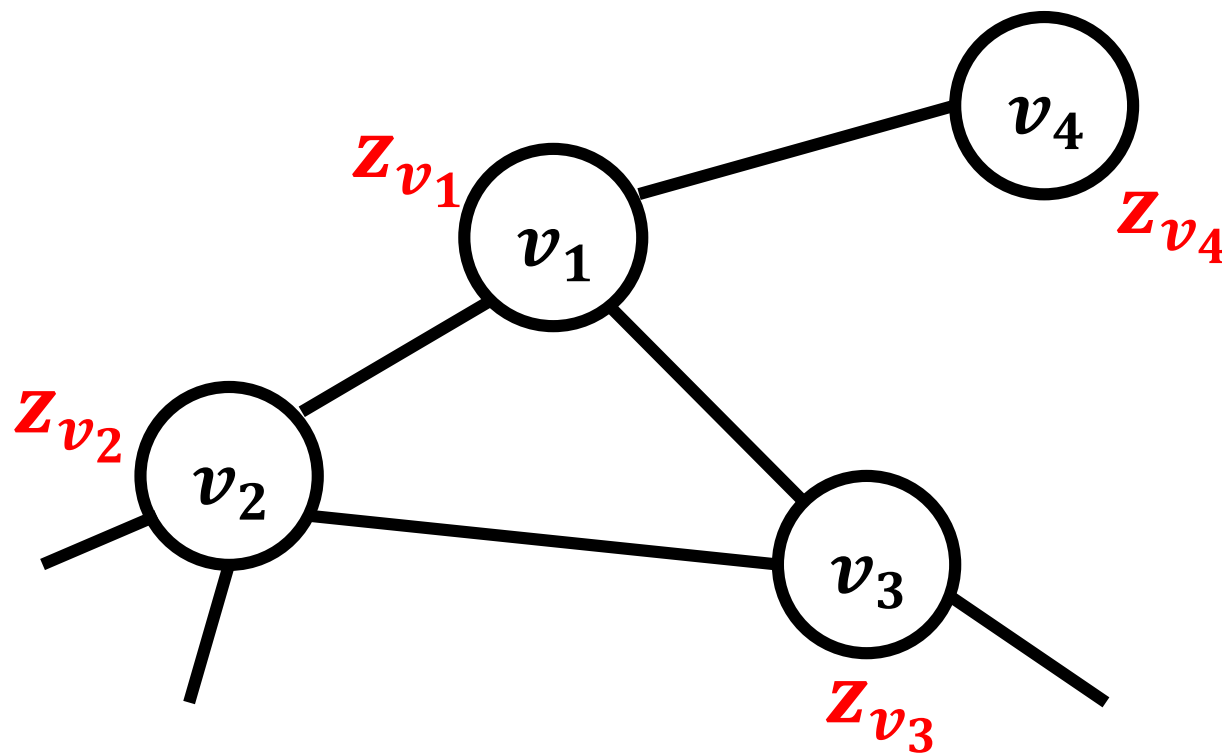
Combine its previous embedding from the last iteration with its aggregated neighborhood vector.



Convolutional Neighborhood Aggregation Method

Step 4

We finally get the vector representation for each node.



Convolutional Neighborhood Aggregation Method

Algorithm 1: Neighborhood-aggregation encoder algorithm. Adapted from [28].

Input : Graph $\mathcal{G}(\mathcal{V}, \mathcal{E})$; input features $\{\mathbf{x}_v, \forall v \in \mathcal{V}\}$; depth K ; weight matrices $\{\mathbf{W}^k, \forall k \in [1, K]\}$; non-linearity σ ; differentiable aggregator functions $\{\text{AGGREGATE}_k, \forall k \in [1, K]\}$; neighborhood function $\mathcal{N} : v \rightarrow 2^{\mathcal{V}}$

Output: Vector representations \mathbf{z}_v for all $v \in \mathcal{V}$

```
1  $\mathbf{h}_v^0 \leftarrow \mathbf{x}_v, \forall v \in \mathcal{V}$ ;
2 for  $k = 1 \dots K$  do
3   for  $v \in \mathcal{V}$  do
4     2  $\mathbf{h}_{\mathcal{N}(v)}^k \leftarrow \text{AGGREGATE}_k(\{\mathbf{h}_u^{k-1}, \forall u \in \mathcal{N}(v)\})$ ;
5     3  $\mathbf{h}_v^k \leftarrow \sigma(\mathbf{W}^k \cdot \text{COMBINE}(\mathbf{h}_v^{k-1}, \mathbf{h}_{\mathcal{N}(v)}^k))$ ;
6   end
7    $\mathbf{h}_v^k \leftarrow \text{NORMALIZE}(\mathbf{h}_v^k), \forall v \in \mathcal{V}$ 
8 end
9 4  $\mathbf{z}_v \leftarrow \mathbf{h}_v^K, \forall v \in \mathcal{V}$ 
```

trainable

Aggregation

Concatenation

Neighborhood Aggregation

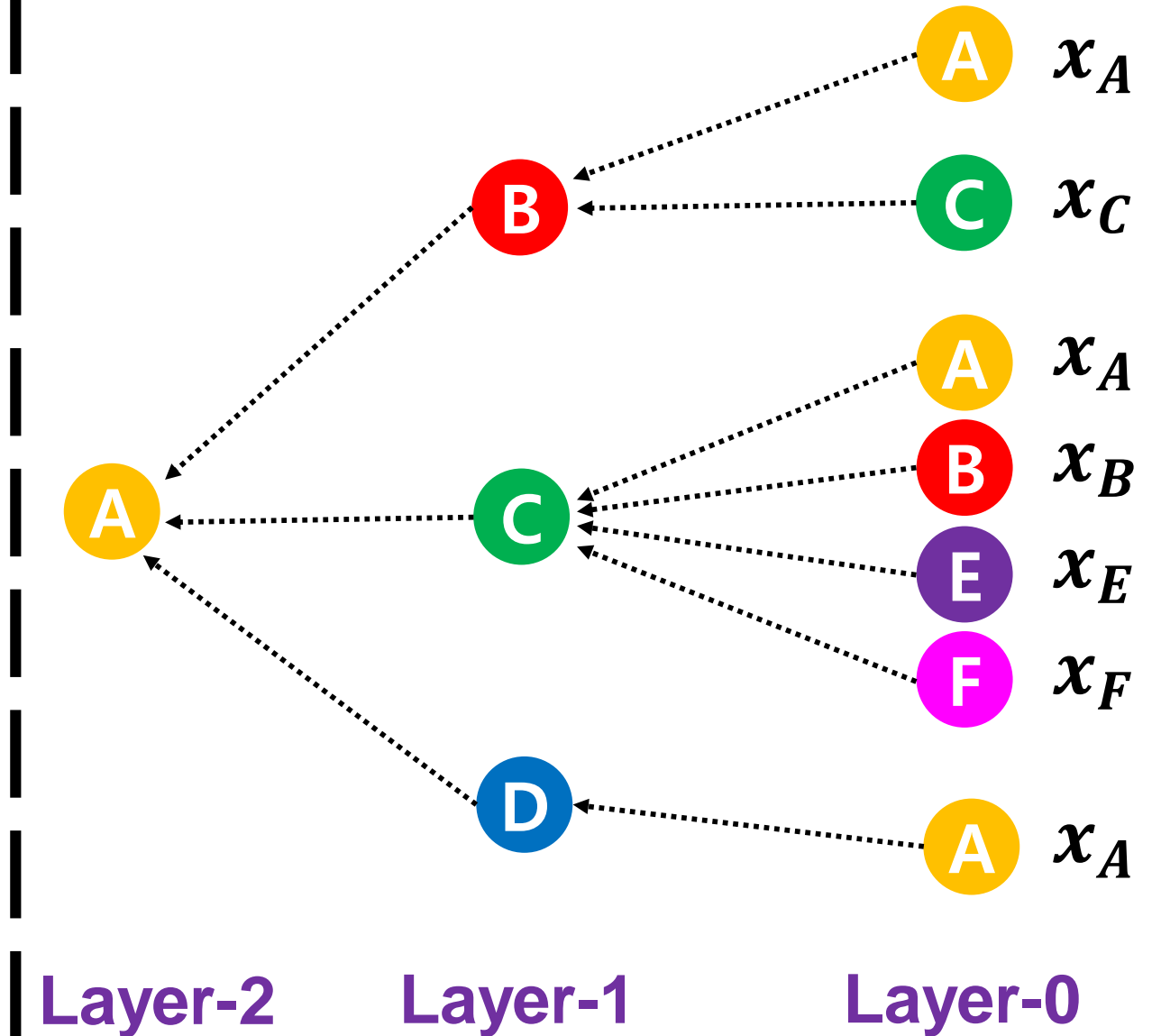
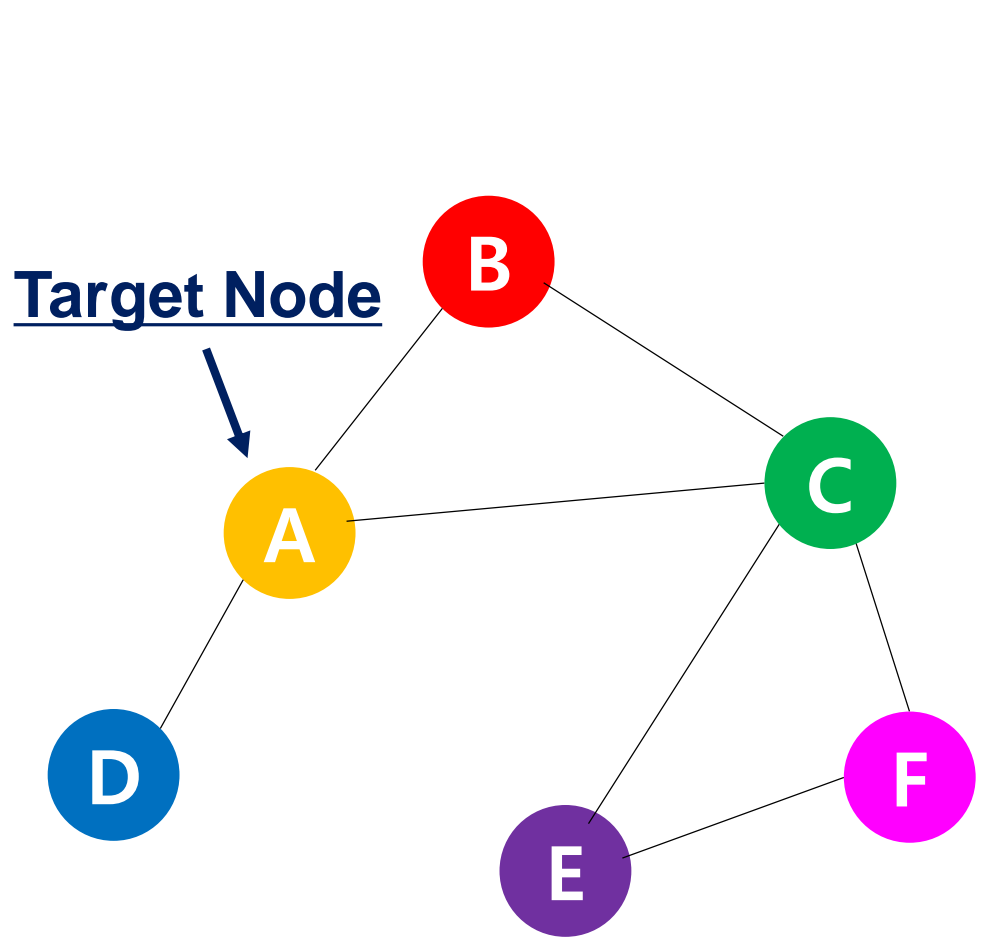


Table of Contents

I. Convolutional Neighborhood Aggregation Method

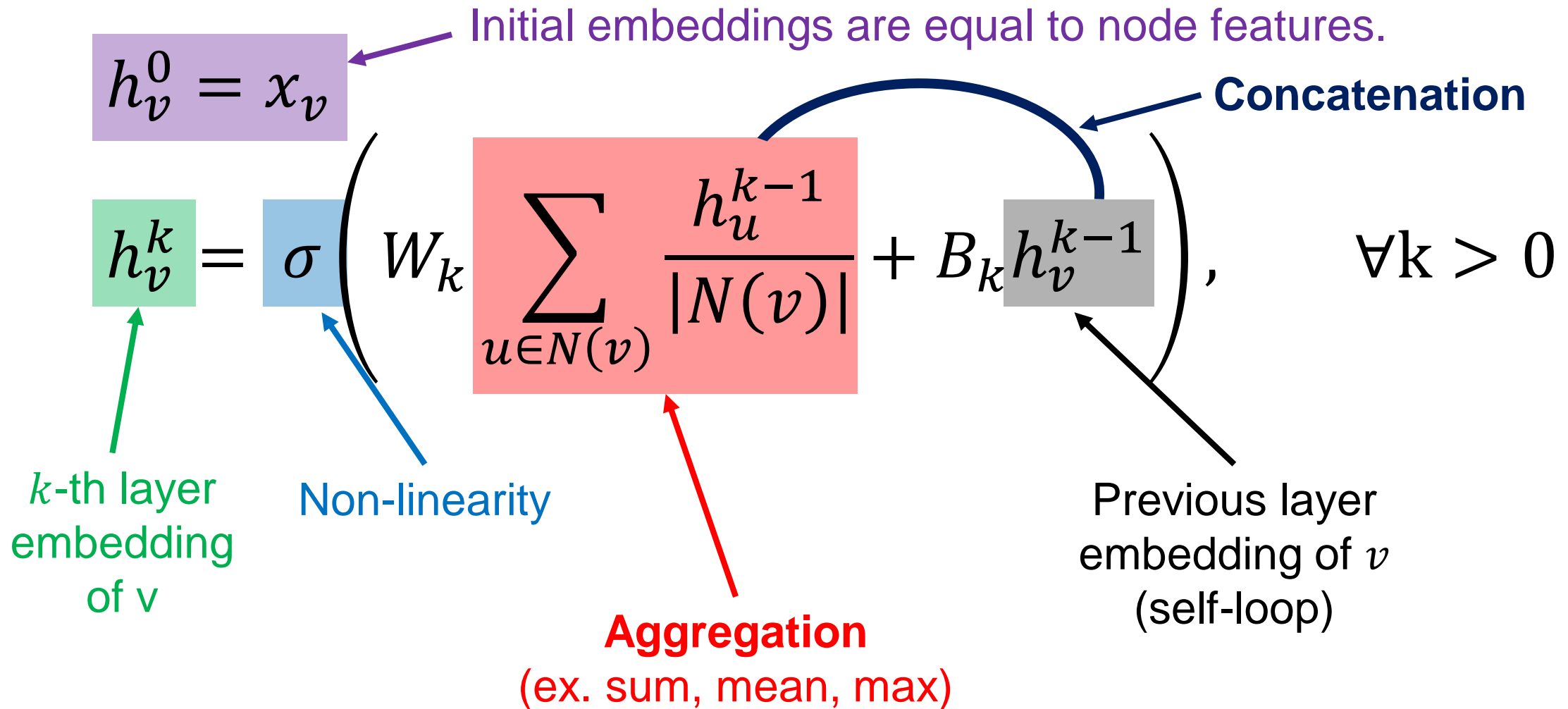
I. Graph Neural Network

II. Graph Convolutional Network

III. GraphSAGE

IV. Graph Attention Network

Graph Neural Network

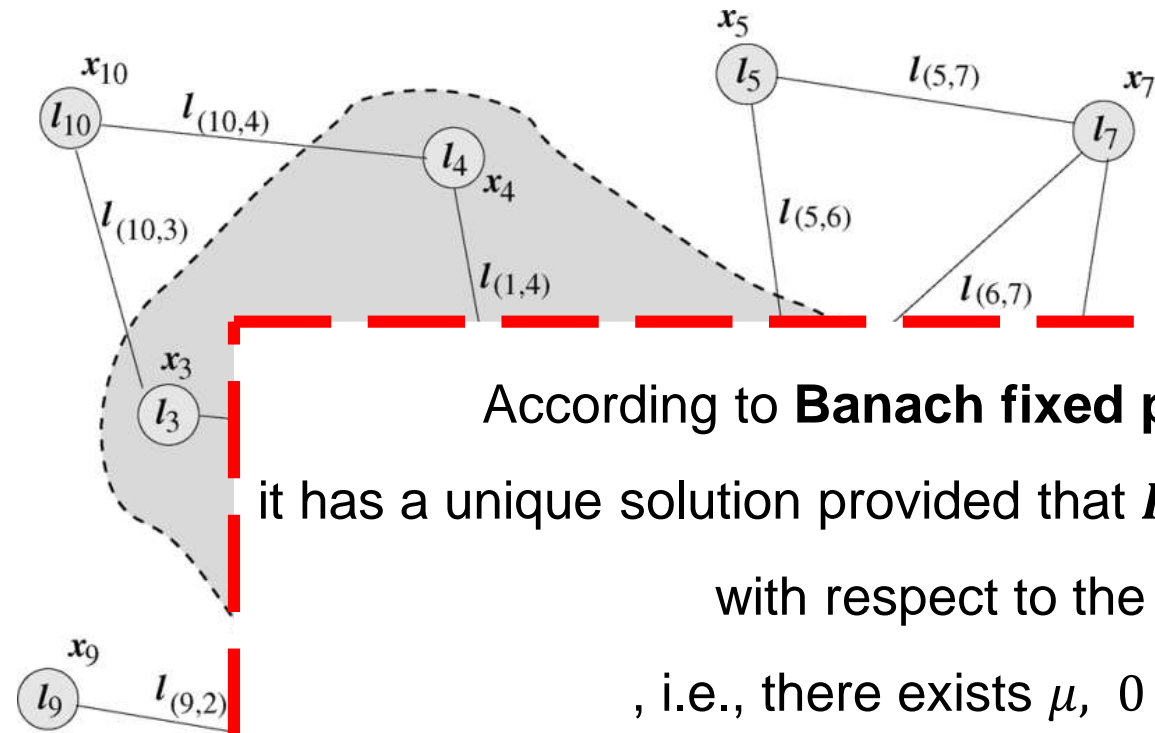


Graph Neural Network

Trainable matrices
(i.e., what we learn)

$$h_v^0 = x_v$$
$$h_v^k = \sigma \left(\boxed{W_k} \sum_{u \in N(v)} \frac{h_u^{k-1}}{|N(v)|} + \boxed{B_k} h_v^{k-1} \right), \quad \forall k \in \{1, \dots, K\}$$
$$z_v = h_v^K$$

- After **K**-layers of neighborhood aggregation, we get output embeddings for each node.
- We can feed these embeddings into any loss function and run stochastic gradient descent to **train the aggregation parameters**.



According to **Banach fixed point theorem**,
it has a unique solution provided that F_w is a **contraction map**
with respect to the state

, i.e., there exists μ , $0 \leq \mu < 1$,

such that $\|F_w(x, l) - F_w(y, l)\| \leq \mu \|x - y\|$ holds for any x, y .

$$x_1 = f_w(l_1, l_{(1,2)}, l_{(1,3)}, x_2, x_3, x_4, x_6, l_2, l_3, l_4, l_6)$$

$\underbrace{l_{(1,2)}, l_{(1,3)}, x_2, x_3, x_4, x_6}_{l_{co}[1]} \quad \underbrace{l_2, l_3, l_4, l_6}_{l_{ne}[n]}$

$$\mathbf{x}_n = f_w(\mathbf{l}_n, \mathbf{l}_{co}[n], \mathbf{x}_{ne}[n], \mathbf{l}_{ne}[n])$$

$$\mathbf{o}_n = g_w(\mathbf{x}_n, \mathbf{l}_n)$$

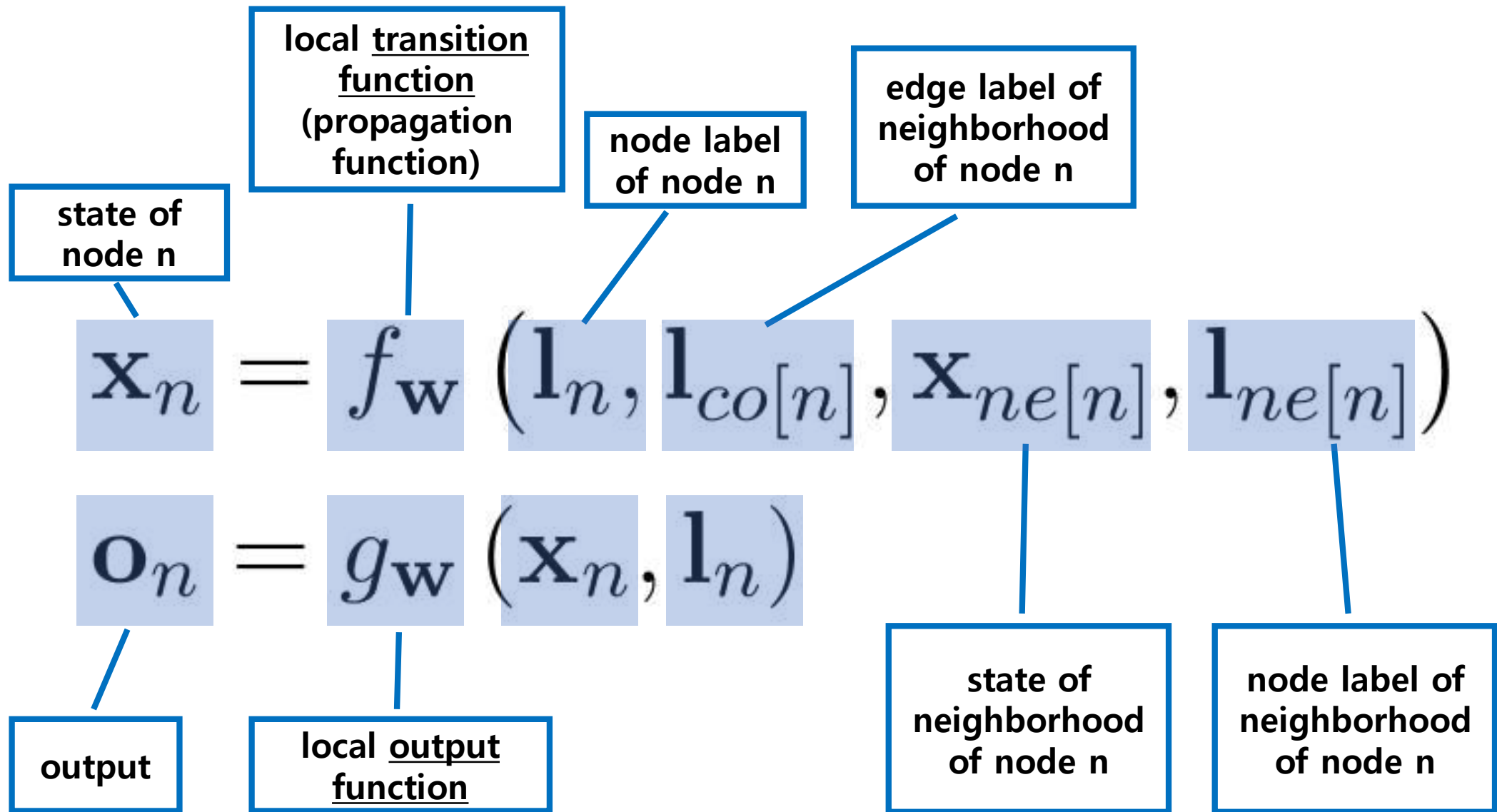


Table of Contents

I. Convolutional Neighborhood Aggregation Method

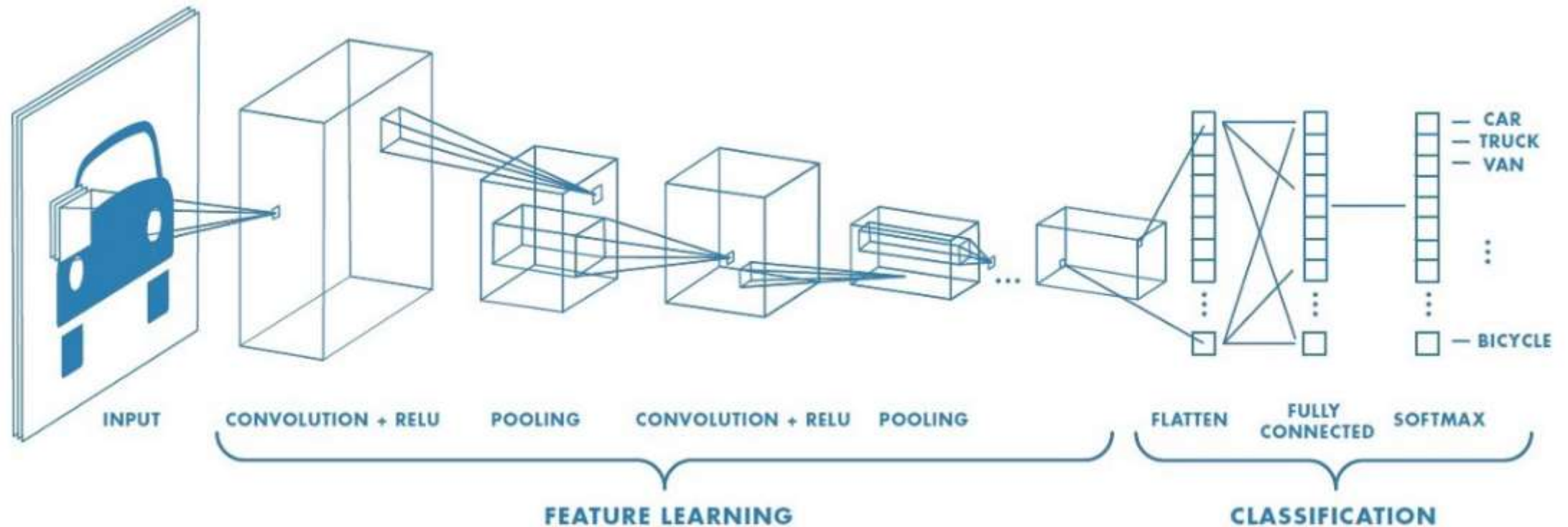
I. Graph Neural Network

II. Graph Convolutional Network

III. GraphSAGE

IV. Graph Attention Network

Graph Convolutional Network



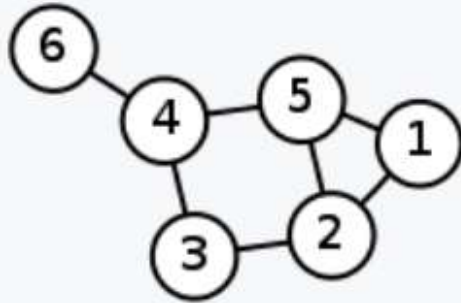
Source : <https://towardsdatascience.com/a-comprehensive-guide-to-convolutional-neural-networks-the-eli5-way-3bd2b1164a53>

Propagation Rule

$$h_i^{(l+1)} = \sigma(W^{(l)} \sum_{j=1}^N h_j^{(l)} + b^{(l)})$$

$$H^{(l+1)} = \sigma(W^{(l)} H^{(l)} + b^{(l)})$$

$$H^{(l+1)} = \sigma(\textcolor{red}{A} W^{(l)} H^{(l)} + b^{(l)})$$

Labelled graph	Adjacency matrix
	$\begin{pmatrix} 0 & 1 & 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \end{pmatrix}$

Graph Convolutional Network

$$H^{(l+1)} = \sigma(\mathbf{A}W^{(l)}H^{(l)} + b^{(l)})$$

$$H^{(l+1)} = \sigma(\tilde{\mathbf{A}}W^{(l)}H^{(l)} + b^{(l)})$$

$$H^{(l+1)} = \sigma(\tilde{\mathbf{D}}^{-\frac{1}{2}}\tilde{\mathbf{A}}\tilde{\mathbf{D}}^{-\frac{1}{2}}W^{(l)}H^{(l)} + b^{(l)})$$

$\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{I}$ (self-loop)

normalization

$$\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{I} \text{ (self-loop)}$$

Table of Contents

I. Convolutional Neighborhood Aggregation Method

I. Graph Neural Network

II. Graph Convolutional Network

III. GraphSAGE

IV. Graph Attention Network

GraphSAGE


- Simple Neighborhood Aggregation:

$$h_v^k = \sigma \left(W_k \sum_{u \in N(v)} \frac{h_u^{k-1}}{|N(v)|} + B_k h_v^{k-1} \right)$$

- GraphSAGE

$$h_v^k = \sigma \left([W_k \cdot \text{AGG}(\{h_u^{k-1}, \forall u \in N(v)\}), B_k h_v^{k-1}] \right)$$

GraphSAGE

- Mean aggregator:  Neighborhood의 정보를 모으는 함수

$$AGG = \sum_{u \in N(v)} \frac{h_u^{k-1}}{|N(v)|}$$

- LSTM aggregator:

$$AGG = LSTM([h_u^{k-1}, \forall u \in \pi(N(v))])$$

- Pooling aggregator:

$$AGG = \gamma(\{Qh_u^{k-1}, \forall u \in N(v)\})$$



Symmetric vector function
(element-wise mean/max)

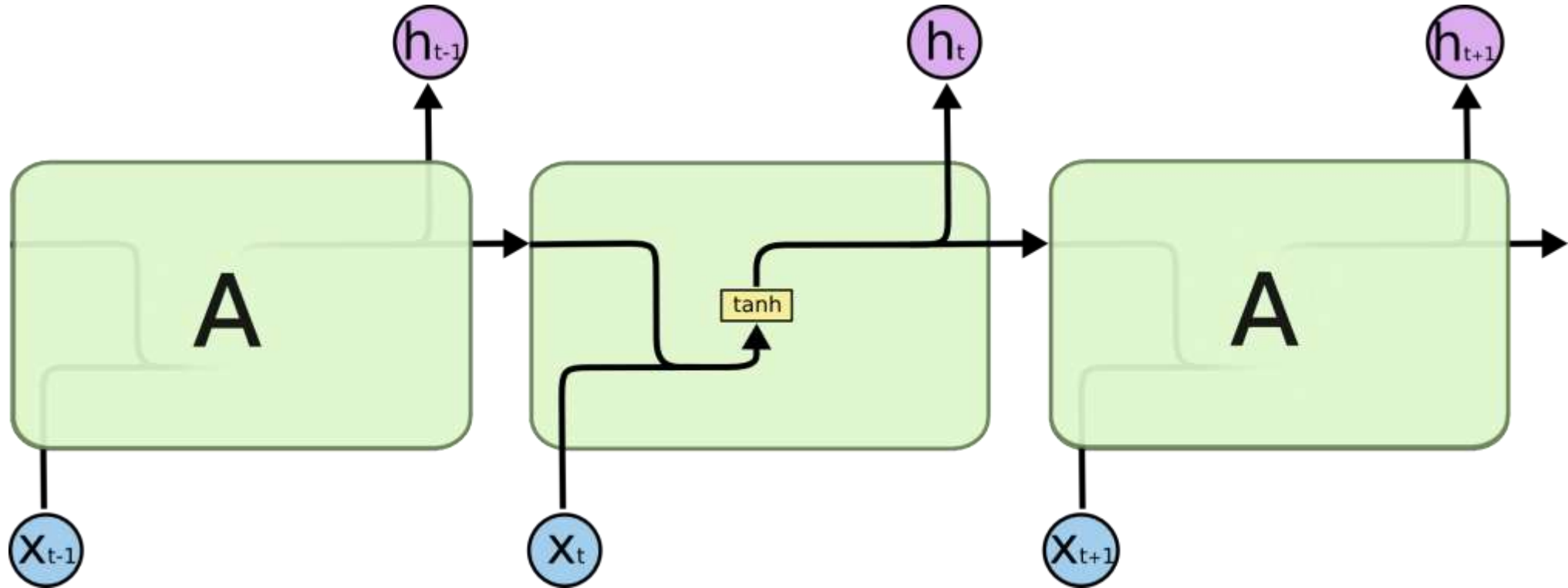
LSTM (Long Short-Term Memory)

- LSTM aggregator:

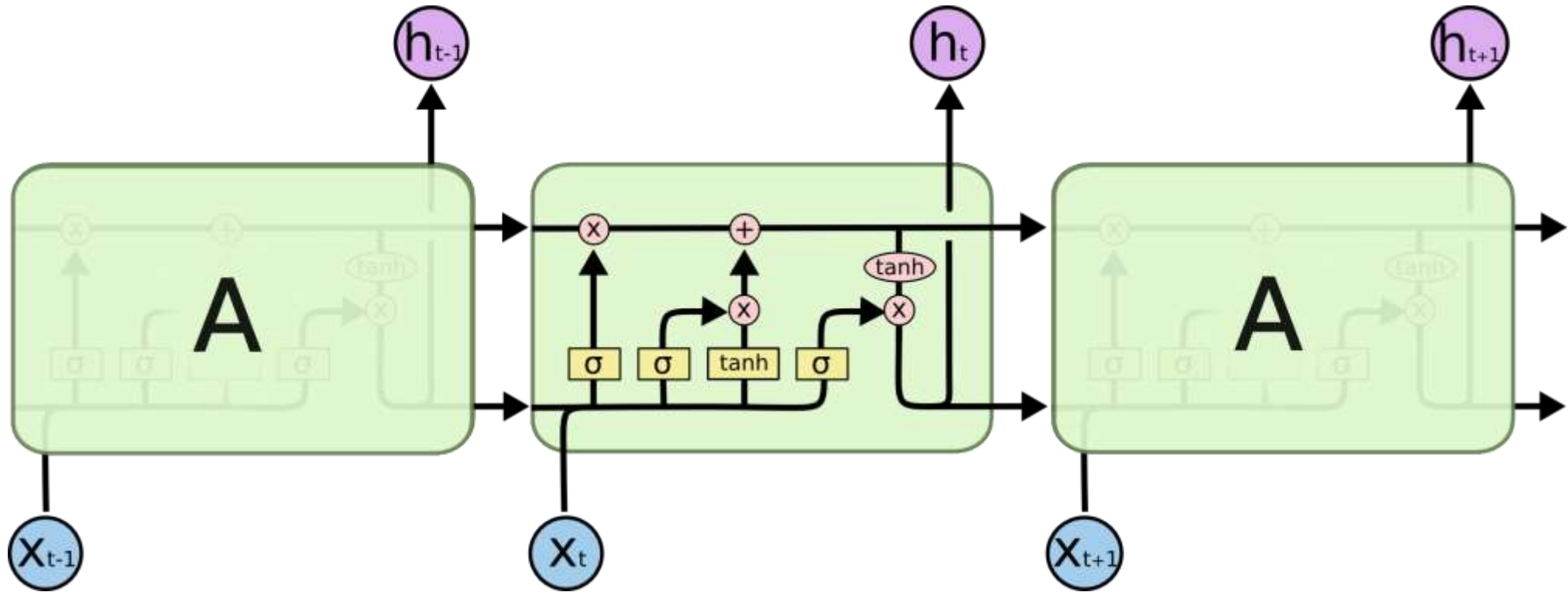
$$AGG = LSTM([h_u^{k-1}, \forall u \in \pi(N(v))])$$

<https://wikidocs.net/152773>

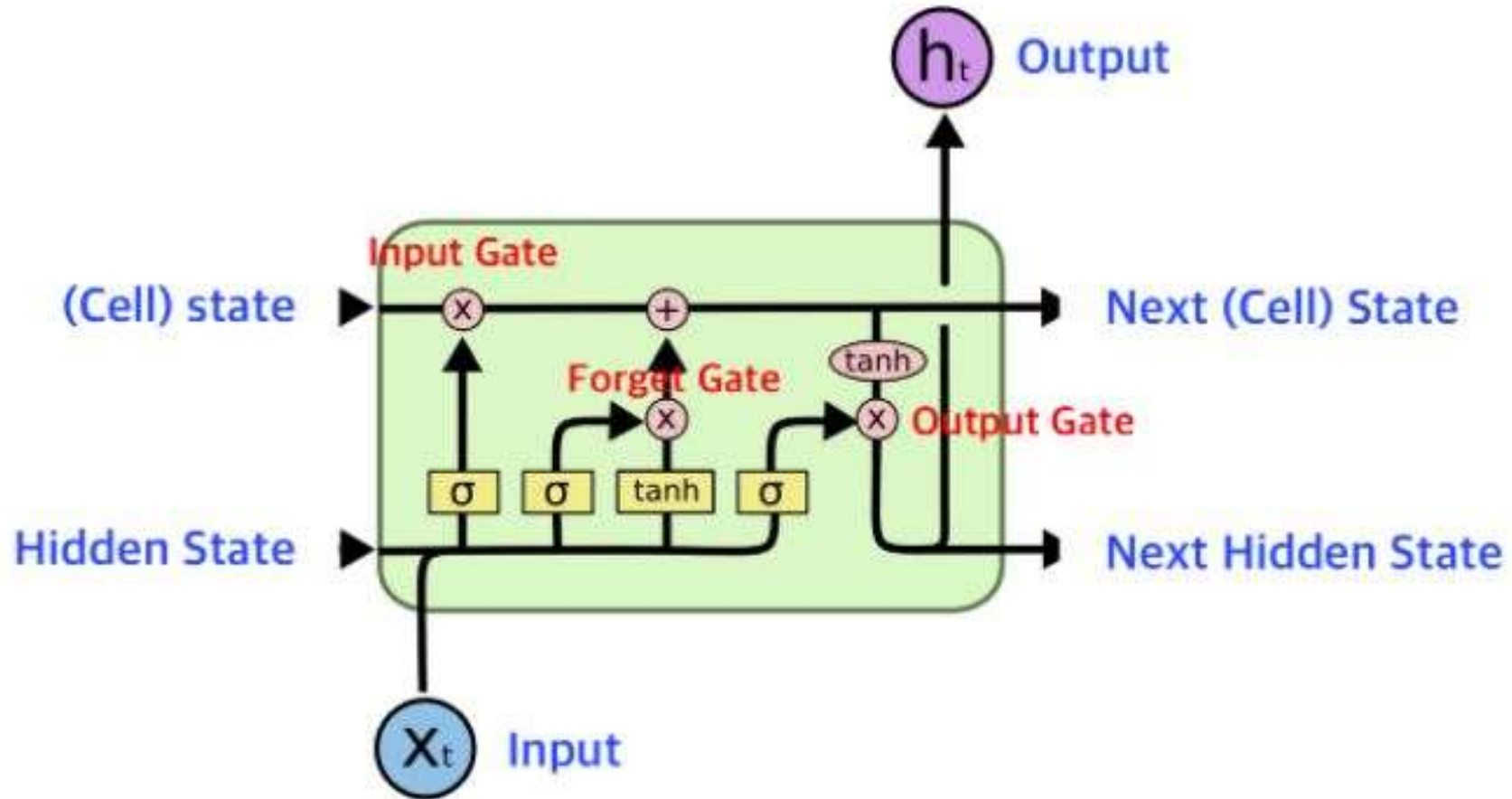
LSTM (Long Short-Term Memory)



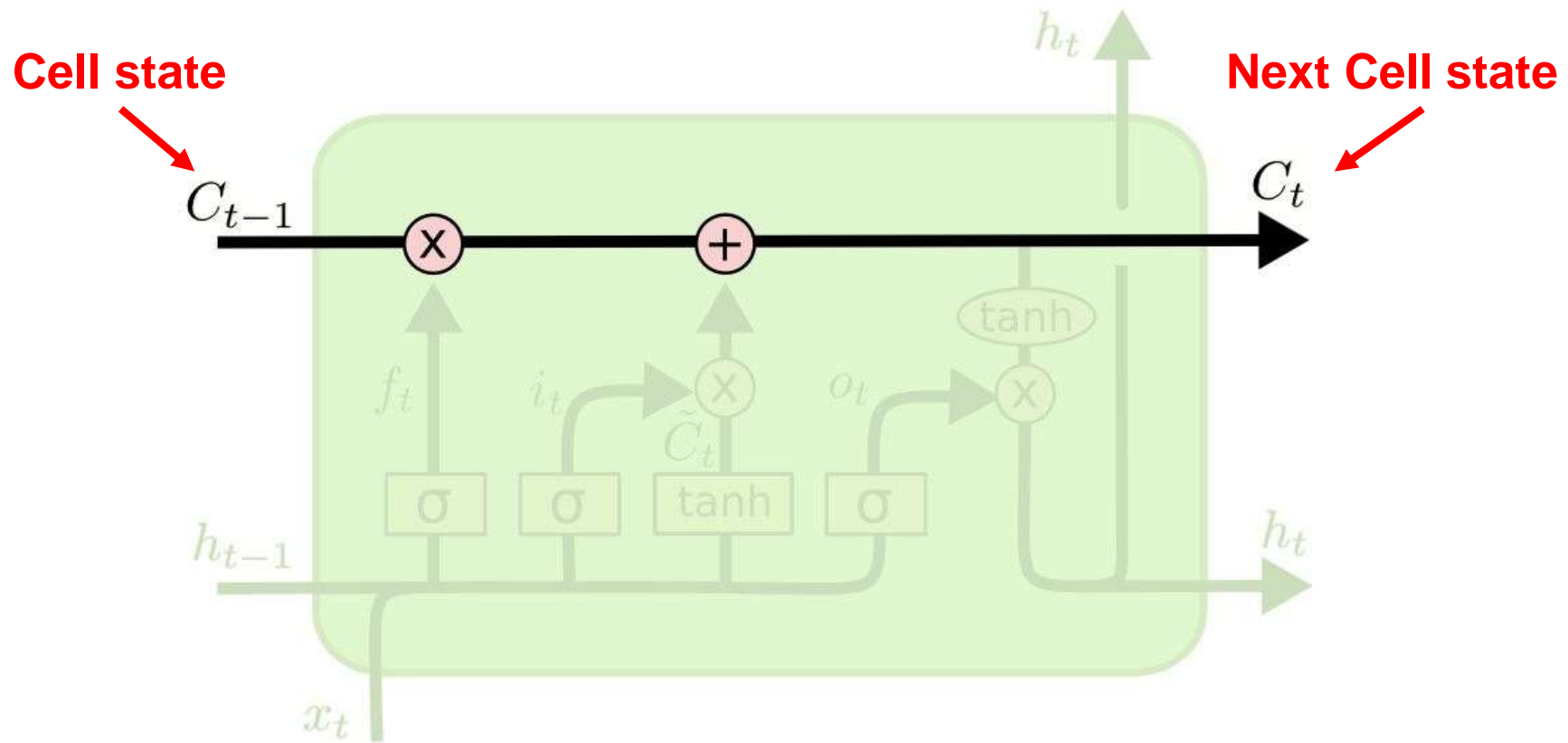
LSTM (Long Short-Term Memory)



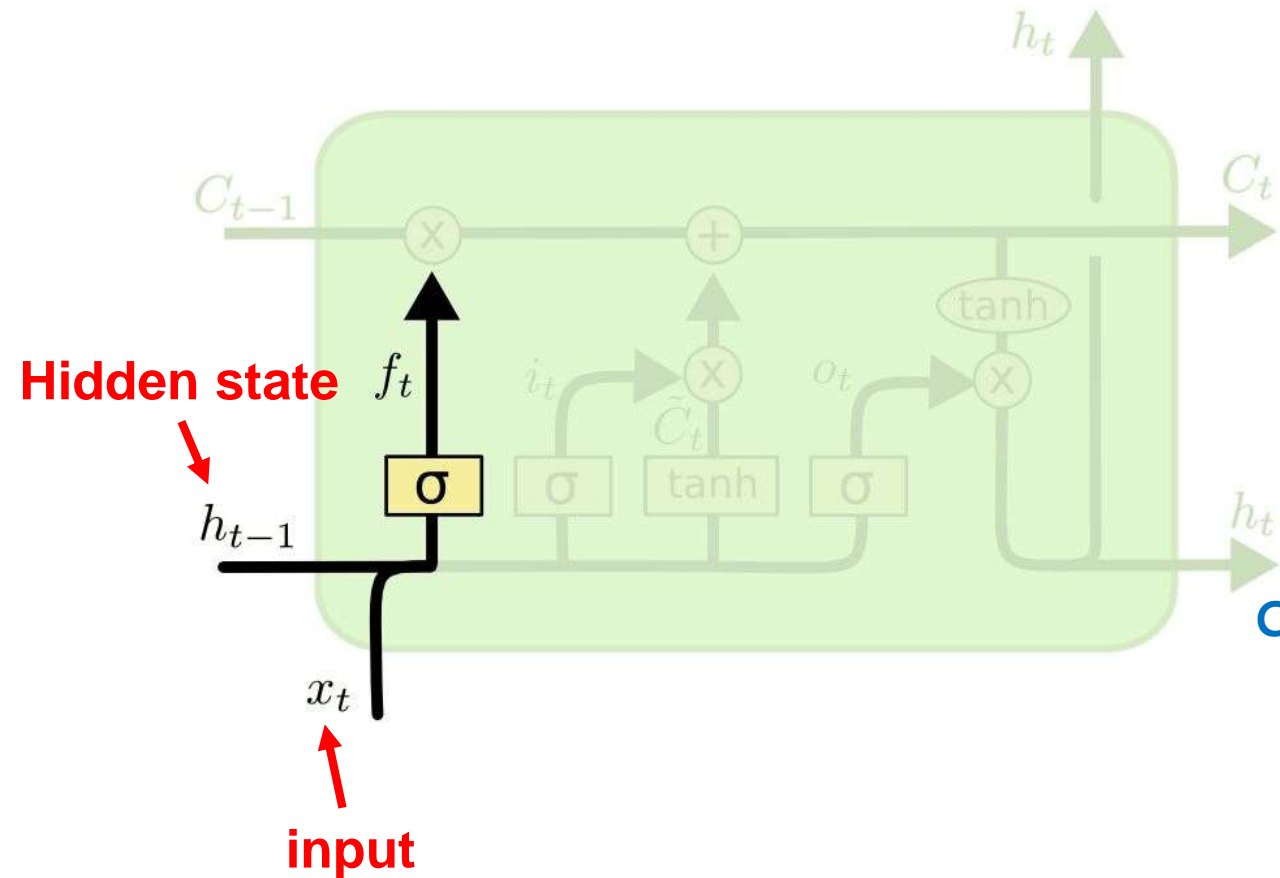
LSTM (Long Short-Term Memory)



Cell State



Forget Gate



- Sigmoid function :

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

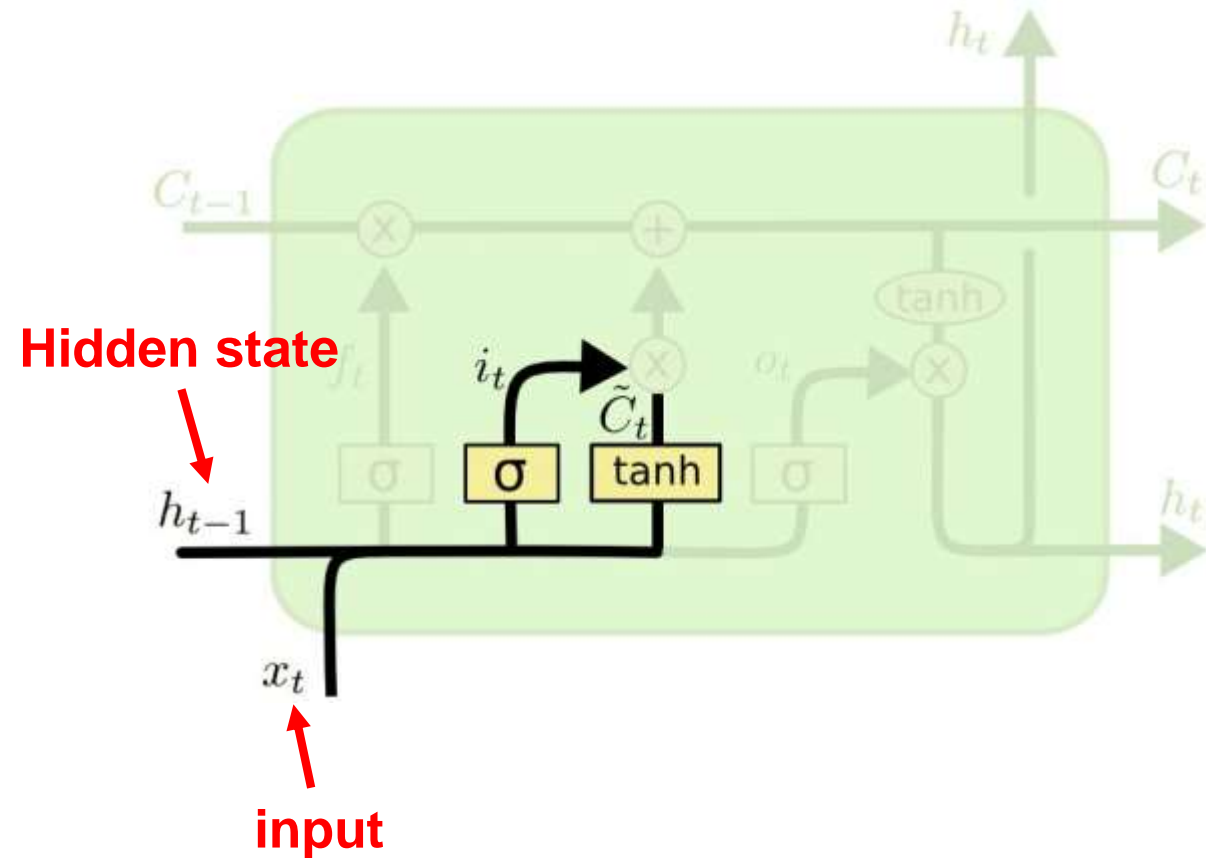
$$\boxed{f_t} = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

Cell state C_{t-1} 을 잊는다



Cell state C_{t-1} 를 기억한다

Input Gate



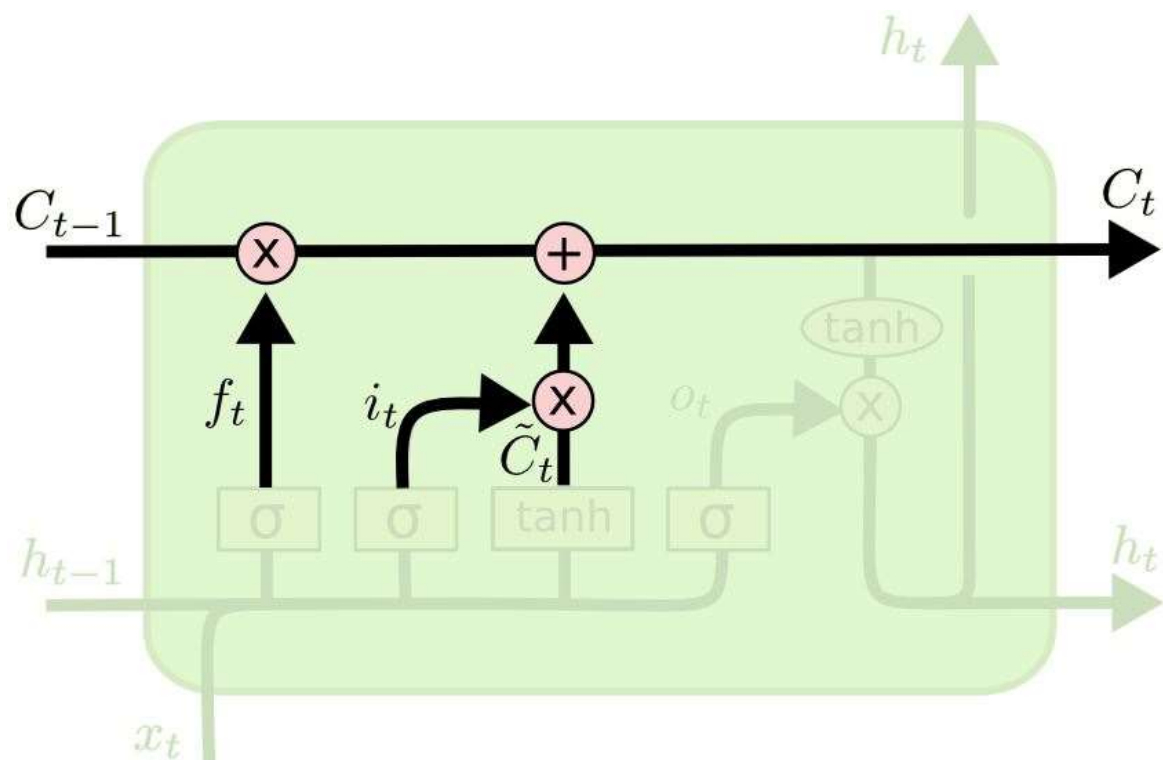
\tilde{C}_t 의 가중치

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

Cell State에 추가될 수 있는 새로운 후보 값

Update

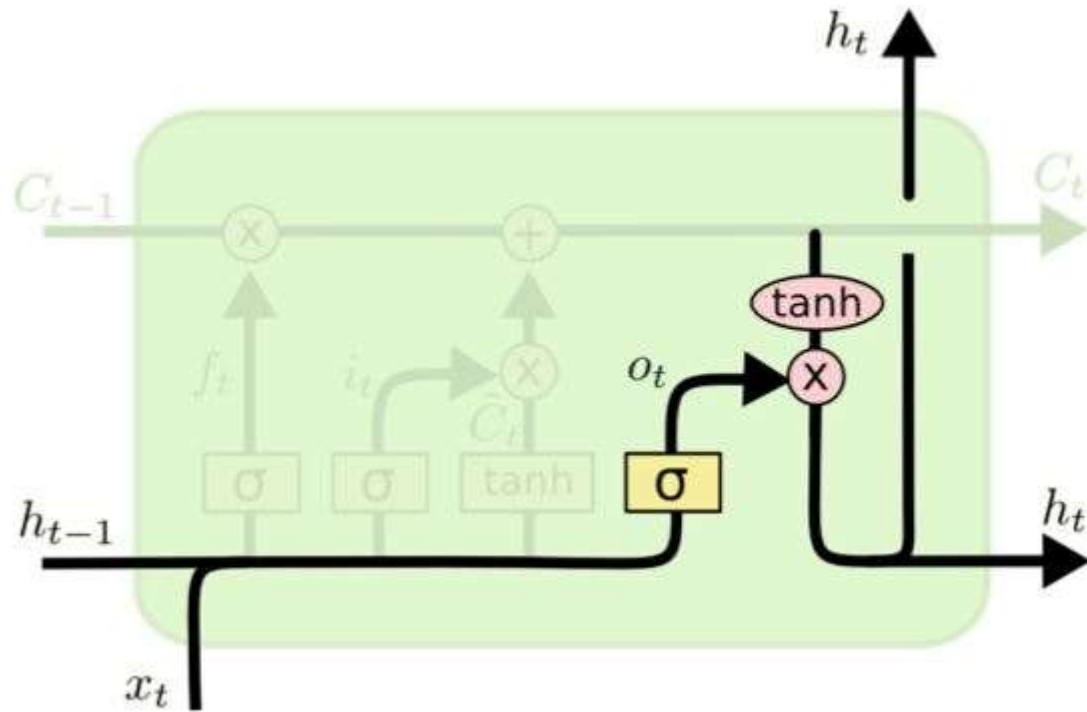


$$C_t = f_t * C_{t-1} + \underbrace{i_t * \tilde{C}_t}_{\text{Input Gate}}$$

Cell state

Forget Gate

Output Gate



$$o_t = \sigma (W_o [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh (C_t)$$

Pooling

- Pooling aggregator:

$$AGG = \gamma(\{Qh_u^{k-1}, \forall u \in N(v)\})$$

1	1	5	6
5	6	7	8
3	2	1	0
1	2	3	4

Max pooling
with 2×2 filters
and stride 2



6	8
3	4

Pooling

- Pooling aggregator:

$$AGG = \gamma(\{Qh_u^{k-1}, \forall u \in N(v)\})$$

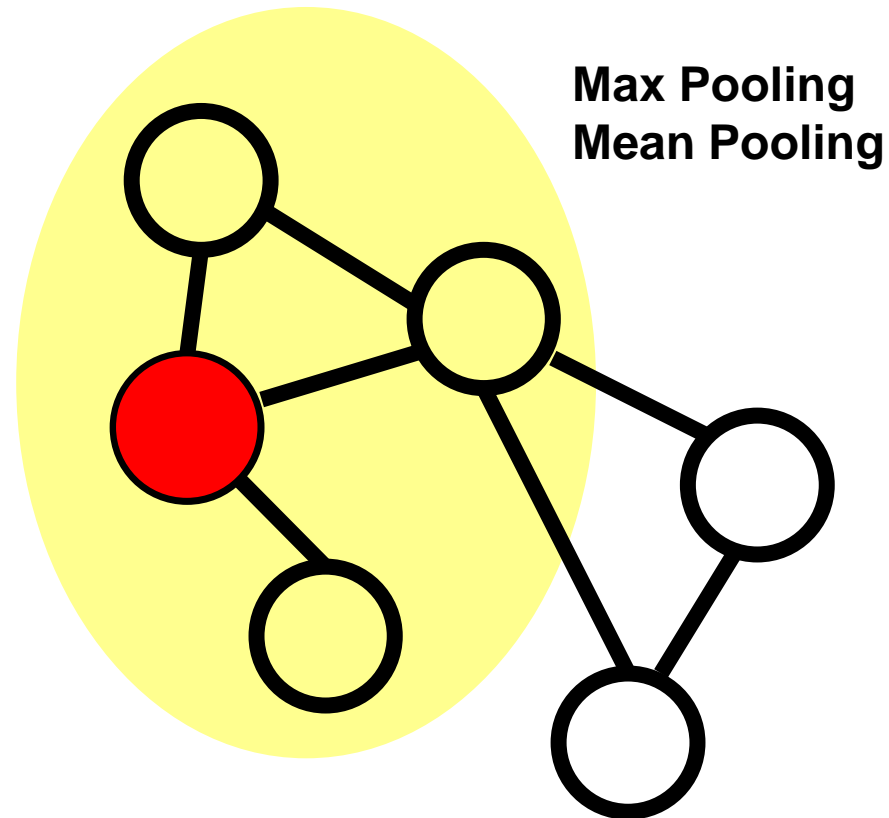


Table of Contents

I. Convolutional Neighborhood Aggregation Method

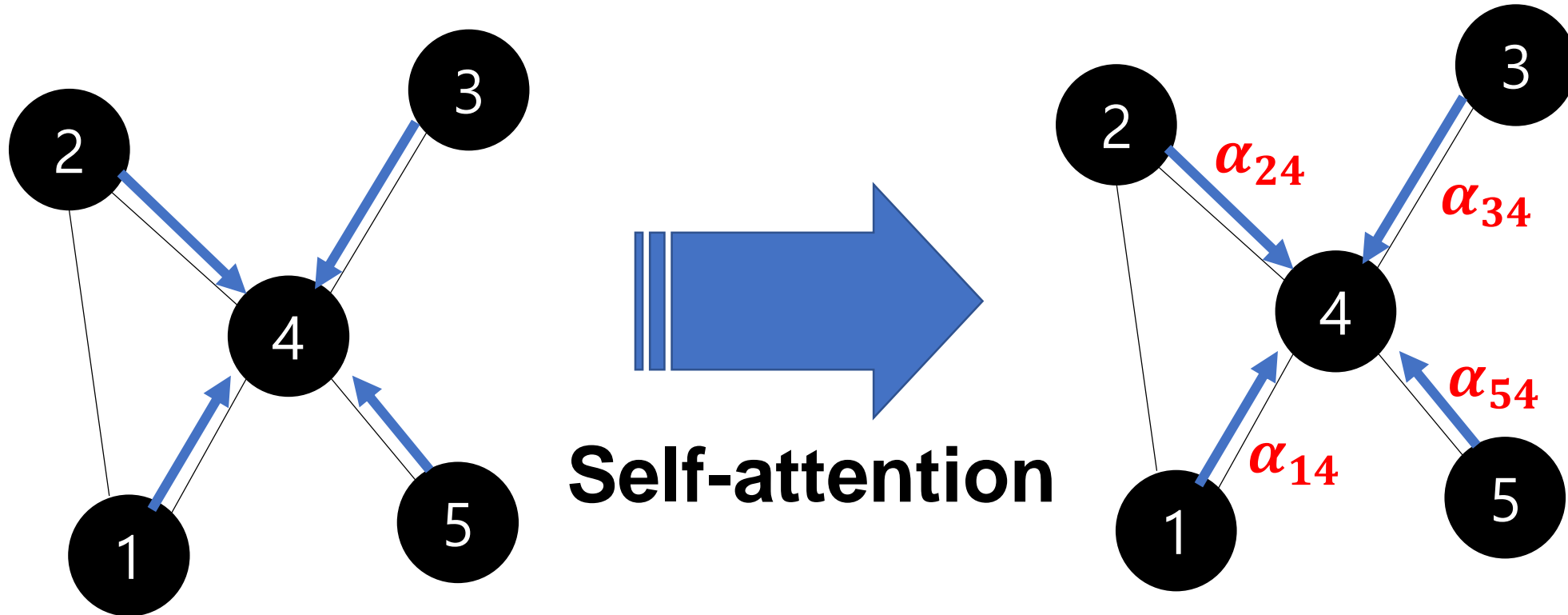
I. Graph Neural Network

II. Graph Convolutional Network

III. GraphSAGE

IV. Graph Attention Network

Graph Attention Network



Graph Attention Network

- We augment basic graph neural network model with attention.

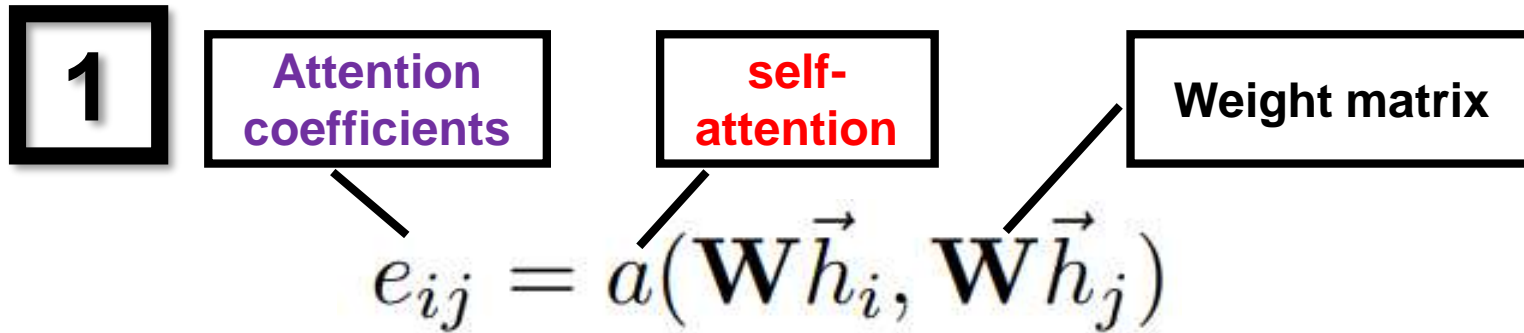
$$h_v^0 = x_v$$

$$h_v^k = \sigma \left(\sum_{u \in N(v)} \alpha_{v,u} W_k \frac{h_u^{k-1}}{|N(v)|} + \alpha_{v,v} B_k h_v^{k-1} \right), \quad \forall k \in \{1, \dots, K\}$$

$$z_v = h_v^K$$

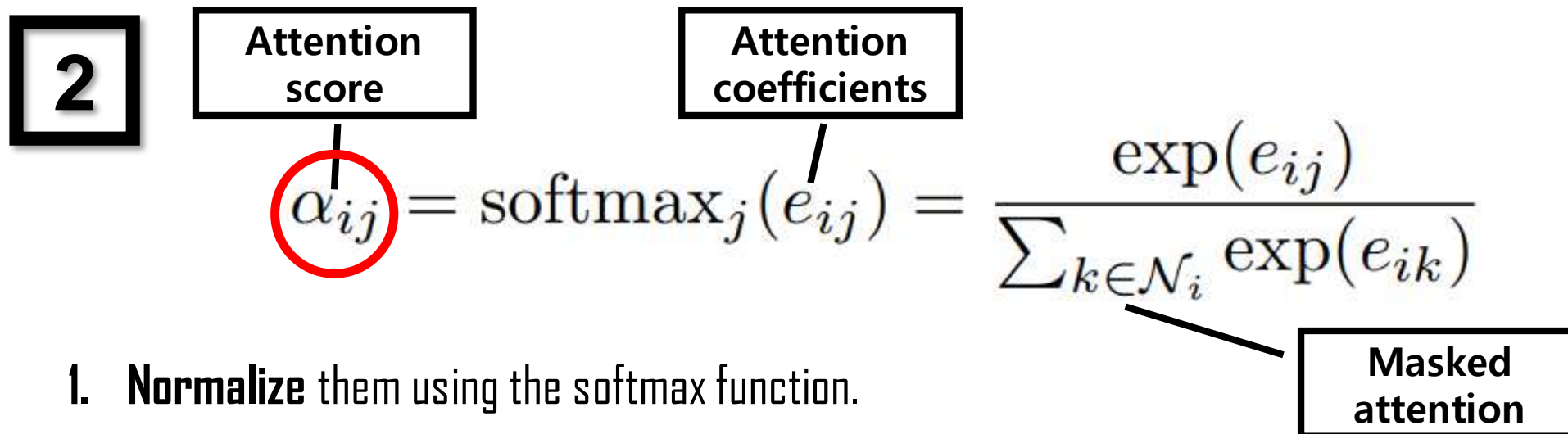
Attention Score





1. Apply a shared linear transformation(\mathbf{W}) to every node.
2. Perform **self-attention**.

Attention coefficients indicate the importance of node j 's features to node i .



1. **Normalize** them using the softmax function.

Self-attention

- **Scaled Dot-Product Attention**

The input consists of queries, keys and values.

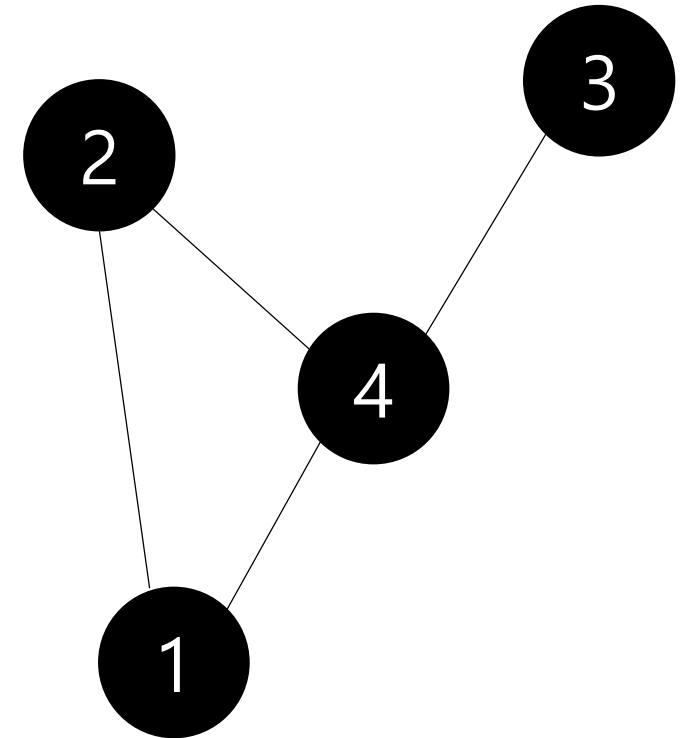
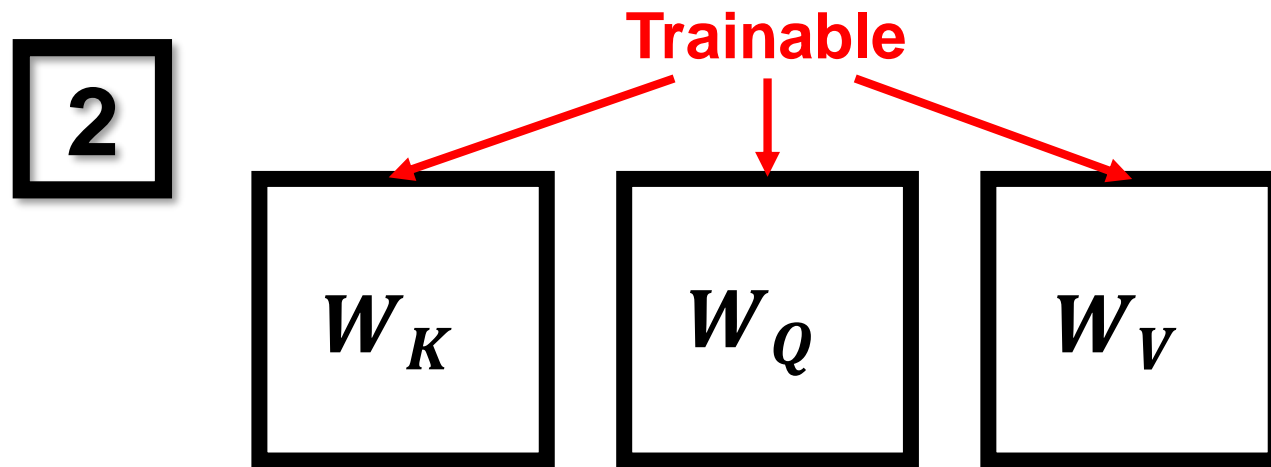
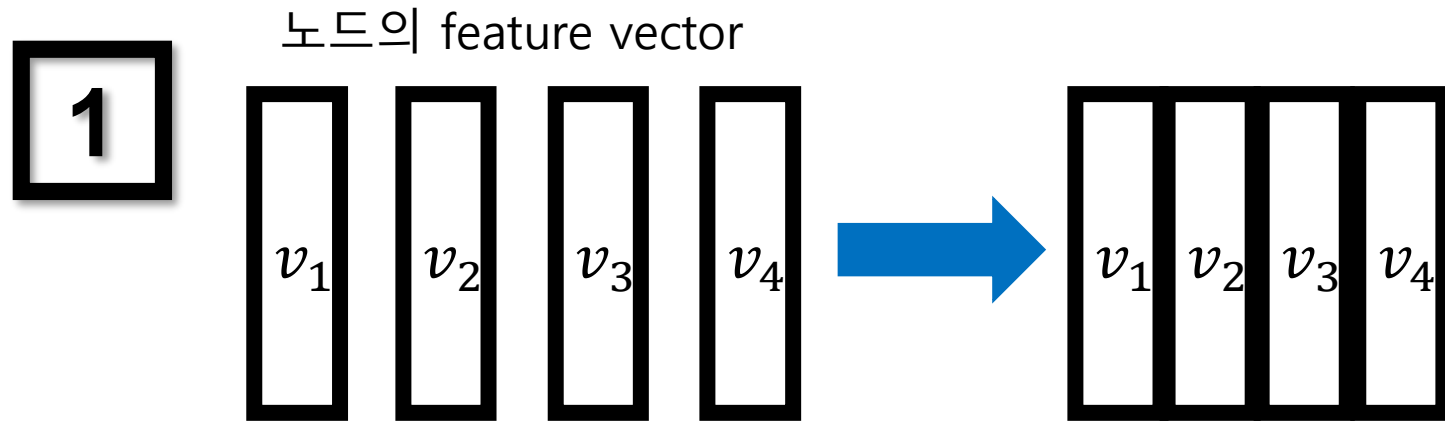
We compute the matrix of outputs as:

Attention Score

$$Att(K, Q, V) = softmax\left(\frac{1}{\sqrt{d_k}} Q K^T\right) V$$

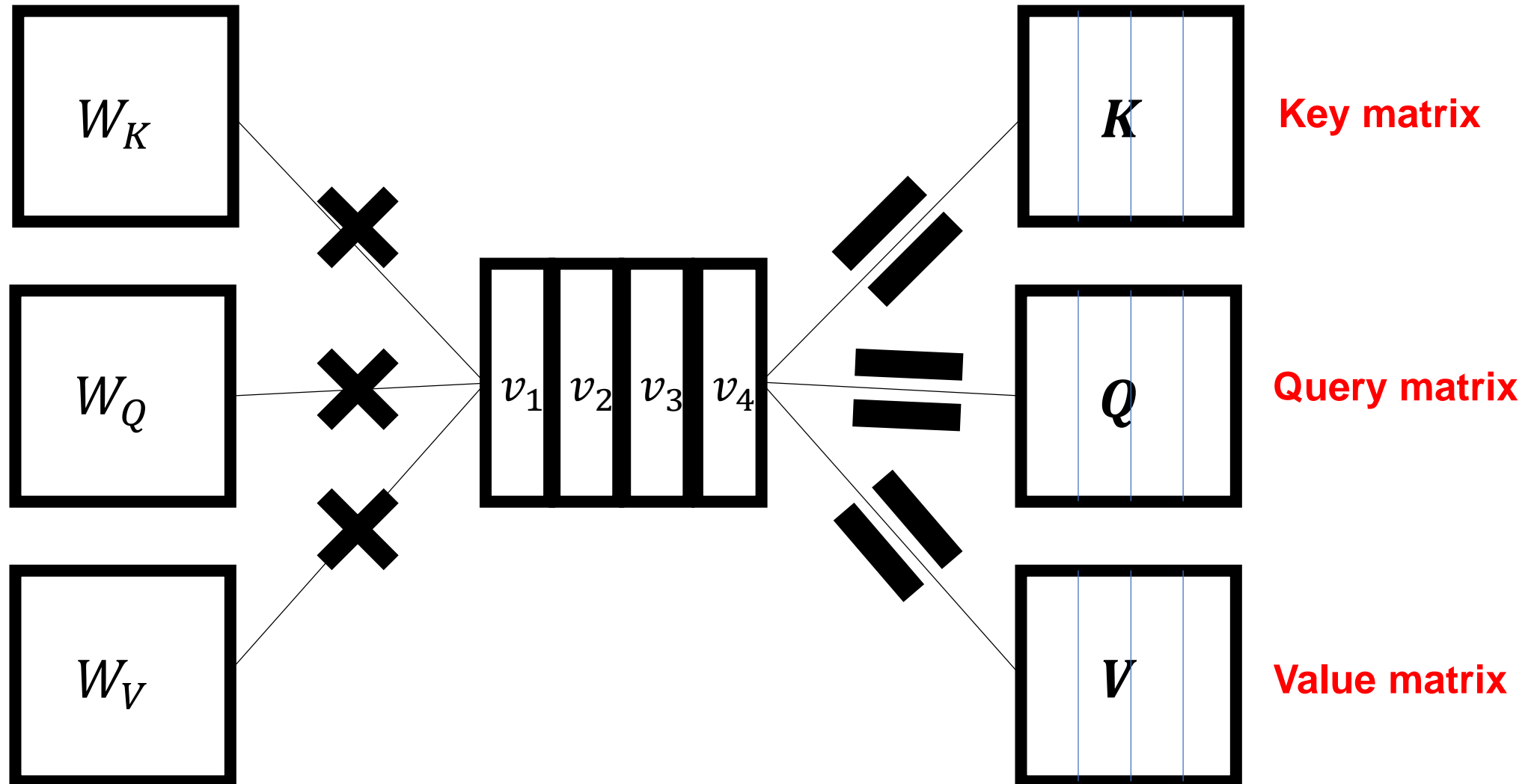
where $K = \begin{pmatrix} k_1^T \\ \vdots \\ k_n^T \end{pmatrix} \in \mathbb{R}^{n \times d_k}$, $Q = \begin{pmatrix} q_1^T \\ \vdots \\ q_m^T \end{pmatrix} \in \mathbb{R}^{m \times d_k}$ and $V = \begin{pmatrix} v_1^T \\ \vdots \\ v_n^T \end{pmatrix} \in \mathbb{R}^{n \times d_v}$.

Self-attention



Self-attention

3

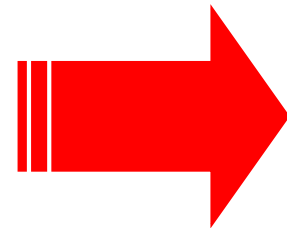


Self-attention

4

The diagram illustrates the calculation of attention coefficients. It shows a matrix multiplication of K^T (a 4x4 matrix with rows labeled v_1, v_2, v_3, v_4) and Q (a 4x4 matrix with columns labeled v_1, v_2, v_3, v_4). The result is a 4x4 matrix, which is then scaled by $\frac{1}{\sqrt{d_k}}$.

$$\begin{matrix} v_1 \\ v_2 \\ v_3 \\ v_4 \end{matrix} \begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \end{bmatrix} K^T \times \begin{matrix} v_1 & v_2 & v_3 & v_4 \\ \begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \end{bmatrix} Q \end{matrix} = \begin{matrix} v_1 & v_2 & v_3 & v_4 \\ \begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \end{bmatrix} \end{matrix} \frac{1}{\sqrt{d_k}}$$



Attention coefficients



Softmax

Attention score

3

$$\vec{h}'_i = \sigma \left(\sum_{j \in \mathcal{N}_i} \alpha_{ij} \mathbf{W} \vec{h}_j \right)$$

1. Compute a linear combination of the features.
2. Apply a *nonlinear* function σ .

expand to
multi-head attention

4

$$\vec{h}'_i = \left\| \begin{matrix} K \\ k=1 \end{matrix} \right\| \sigma \left(\sum_{j \in \mathcal{N}_i} \alpha_{ij}^k \mathbf{W}^k \vec{h}_j \right)$$

K independent attention mechanism

- $\|$: concatenation
- α_{ij}^k : normalized attention coefficients computed by the **k-th** attention mechanism(a^k)
- \mathbf{W}^k : weight matrix of k-th attentional head

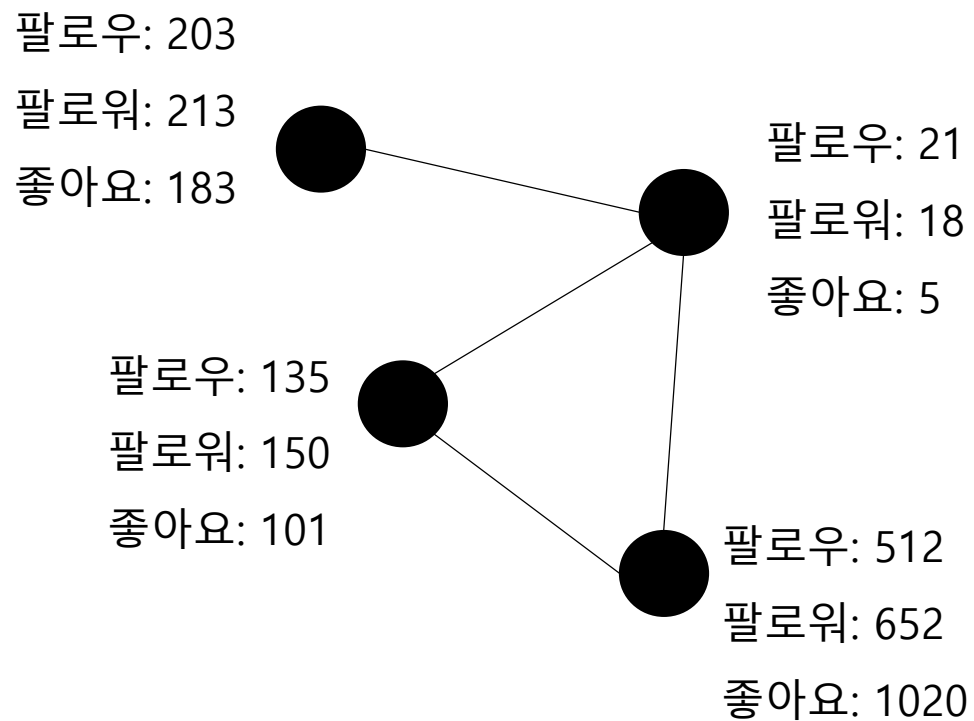
Multi-head Attention

ex) Instagram

팔로워 $\leftarrow head_1$

팔로우 $\leftarrow head_2$

좋아요 $\leftarrow head_3$



Multi-head Attention

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

where $\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$

and the projections $W_i^Q \in \mathbb{R}^{d_{\text{model}} \times d_k}$, $W_i^K \in \mathbb{R}^{d_{\text{model}} \times d_k}$, $W_i^V \in \mathbb{R}^{d_{\text{model}} \times d_v}$ and $W^O \in \mathbb{R}^{hd_v \times d_{\text{model}}}$.

Thank you for listening.