

Graph Convolutional Neural Networks

* Graphs are irregular data structures, typically used to represent hierarchical relationships

* No parent and child relationship exist

$\{ \text{DAG} \Rightarrow \text{Tree} \}$ } with loops & self loops
and directed edges along with weights they become very general & powerful.

$$G = (V, E)$$

Set of Verticies V
 (v_1, v_2, v_3, \dots)

Set of Edges E
 (e_1, e_2, \dots)

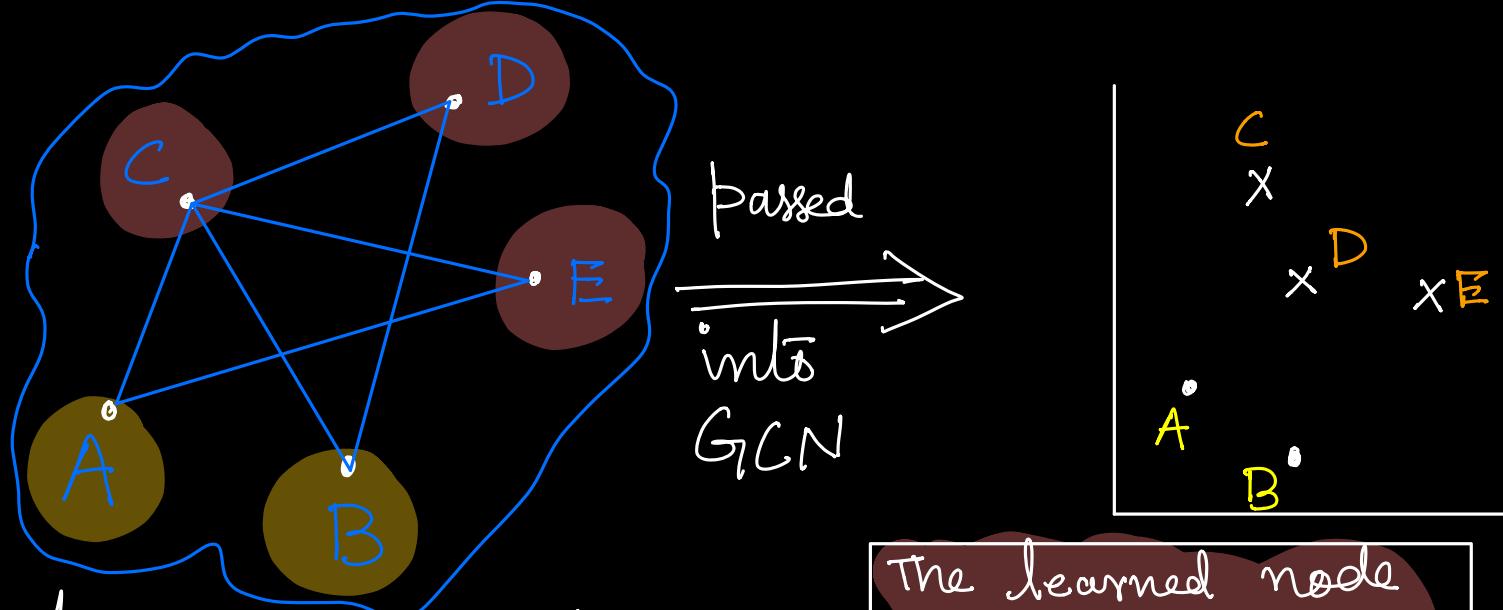
It has been observed that for learning related graph embedded problems

Graph CNN can be very powerful learning tool.

2-types of GCN's
Spectral & Spatial

Several problems can be easily embedded into graphs and can be solved using existing & efficient graph algorithms for traversal, shortest path.

Even with few randomly intiyed GCN layers can extract meaningful node's feature representation.



passed
into
GCN

$C = [x_1, x_2]$
 $D = [x_3, x_4]$
 $E = [x_5, x_6]$

(A, B, C, D, E)

(AC, AE, BC, BD, CD, CE)

Adjacency Matrix \mathbf{A} :

	A	B	C	D	E
A	0	0	1	0	1
B	0	0	1	1	0
C	1	1	0	1	1
D	0	1	1	0	0
E	1	0	1	0	0

	1	2	3	---	F_0
v_1	f_A^1	f_A^2	f_A^3	---	$f_A^{F_0}$
v_2	f_B^1	f_B^2	f_B^3	---	$f_B^{F_0}$
v_3	-	-	-	-	-
v_N	f_E^1	f_E^2	---	---	$f_E^{F_0}$

Each node $\{A, B, C, D, E\}$ have a feature f_A (vector) associated to it. Let us assume that the size of input feature is F^0 . i.e each node have a F^0 dimensional feature vector viz. f_A, f_B, f_C, f_D, f_E .

I/P Feature Matrix (X)

$\Rightarrow N \times F^0$

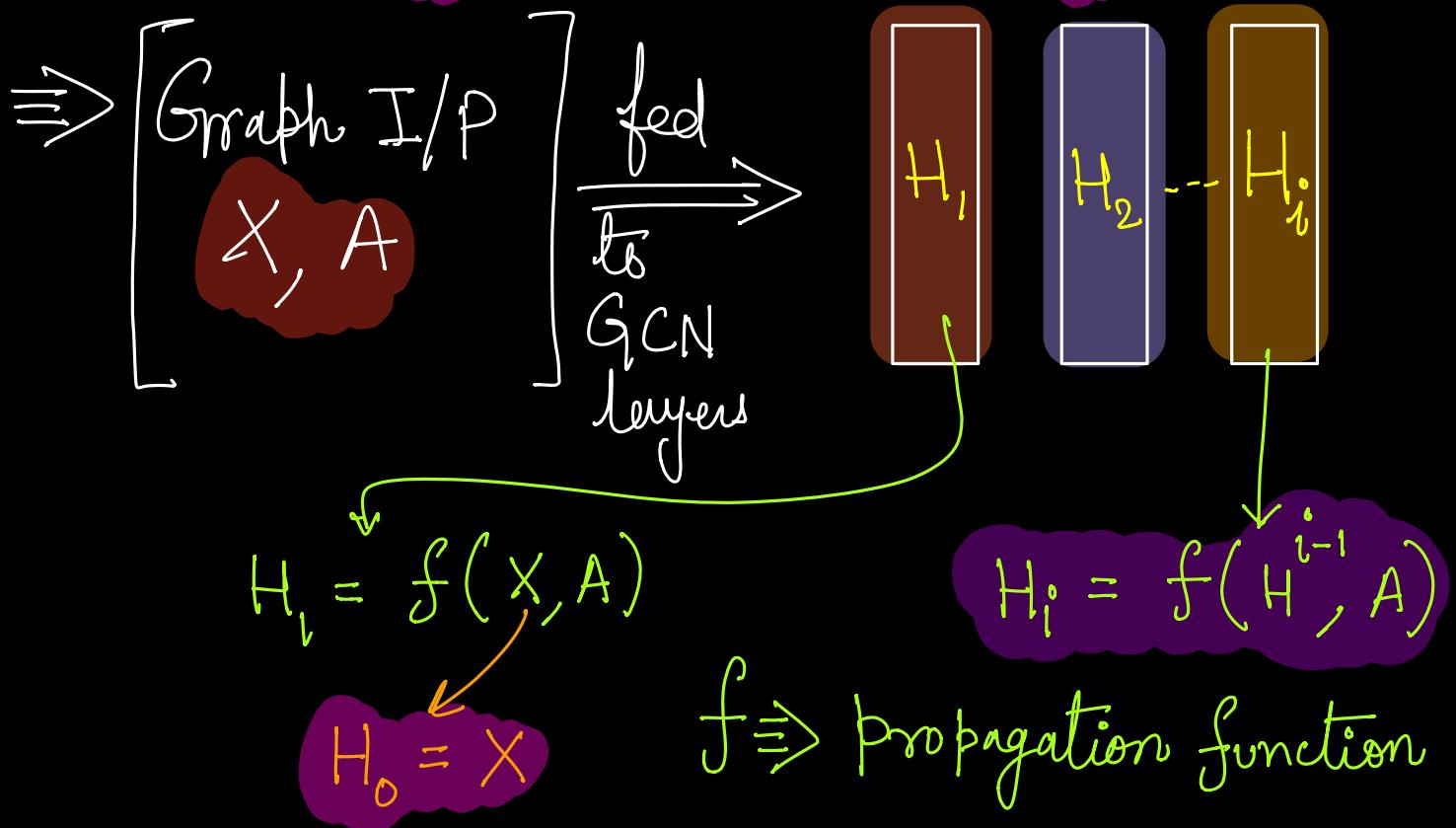
Hence from any given graph $G(V, E)$ we will have 2 matrices extracted as input:

$$|V| = N$$

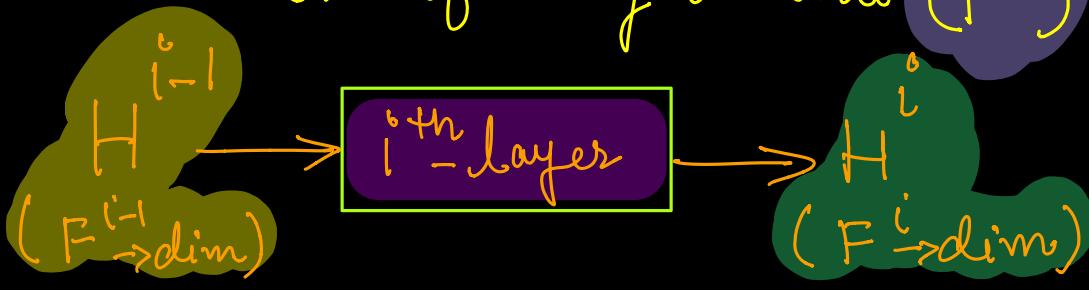
① Adjacency matrix (A) and ② Feature matrix (X)

$$[N \times N]$$

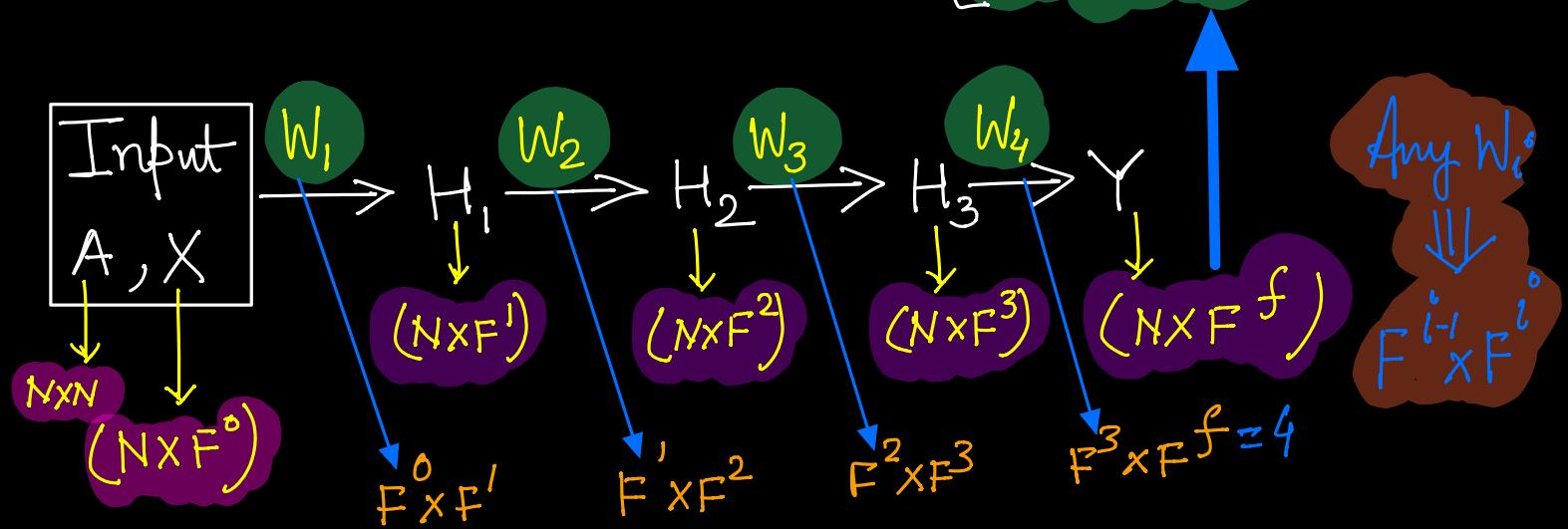
$$[N \times F]$$



① We have the O/P of any (i^{th})-hidden layer (H_i), taking a feature of size (F^{i-1}) corresponding to any node and transforming it into (F^i) dimensions.



Hence the feature matrix after $[i]$ hidden layers will be of size $[N \times F^i]$



- * At each layer the node embedding feature vector got aggregated (utilizing its neighbourhood (A)) to form the next layer features using the learned propagation rule f^i parameterized over W_i

$$\therefore H_i = f_{W_i}(H_{i-1}, A)$$

$$\Rightarrow \text{Say } H_3 = f_{W_3}(H_2, A)$$

H_2 $\in [N \times F^2]$
 $A \in [N \times N]$
 $F^2 \rightarrow f_{W_3} \rightarrow F^3$
 $F^2 \times F^3$

The most simplest rule can be :

$$H_i = f_{W_i}(H_{i-1}, A) = \sigma(A * H_{i-1} * W^i)$$

(N x N)

(N x F^{i-1}) (F^{i-1} x F^i)

(N x F^i)

Talks about node inter-connection and neighbourhood

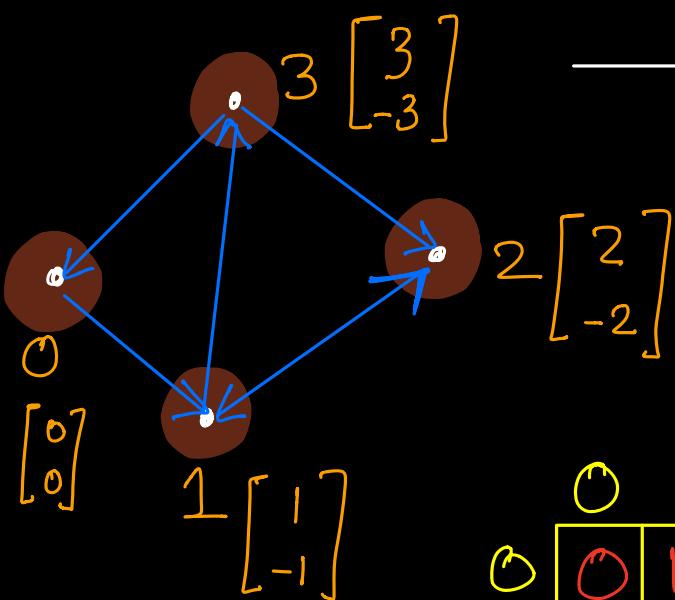
For simplicity let us assume
 σ = linear function

$$\therefore H_i = A H_0 W_0 = A X W_0$$

$$= A X$$

$$\left(\begin{array}{l} W = \frac{1}{\|} \\ \text{initialization} \end{array} \right)$$

Filtering as done in any fully connected network.



Considering such a simple graph with random but easy to use node initial embeddings to understand what is happening.

$$A =$$

	0	1	2	3
0	0	1	0	0
1	0	0	1	1
2	0	1	0	0
3	1	0	1	0

0	0	0
1	1	-1
2	2	-2
3	3	-3

$$\therefore H_1 = f(X, A) = AX = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \end{bmatrix} X \begin{bmatrix} 0 & 0 \\ 1 & -1 \\ 2 & -2 \\ 3 & -3 \end{bmatrix}$$

Any (i^{th}) node feature becomes the addition of features of its neighbours.

$$0 = [0, 0] \Rightarrow [1, -1]$$

$$1 = [1, -1] \Rightarrow [2, -2] + [3, -3] = [5, -5]$$

$$2 = [2, -2] \Rightarrow [1, -1]$$

$$3 = [3, -3] \Rightarrow [0, 0] + [2, -2]$$

$$\begin{bmatrix} 1 & -1 \\ 5 & -5 \\ 1 & -1 \\ 2 & -2 \end{bmatrix}$$

Neighbourhood feature aggregation just as it happens in Convolutional neural network

Problem-01 \rightarrow But since there are no self loops their own features are ignored. In matrix $[A]$ diagonals are "0".

Solution

Just need to add self loop. $[A] \Rightarrow [\hat{A}]$
 $\Rightarrow \hat{A} = A + I$ (Identity matrix)

$$\hat{A} \cdot X = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 0 & 1 & 1 & 1 \\ 0 & 1 & 1 & 0 \\ 1 & 0 & 1 & 1 \end{bmatrix} * \begin{bmatrix} 0 & 0 \\ 1 & -1 \\ 2 & -2 \\ 3 & -3 \end{bmatrix} = \begin{bmatrix} 1 & -1 \\ 6 & -6 \\ 3 & -3 \\ 5 & -5 \end{bmatrix}$$

nodes
own
features are
not ignored.

Problem-02 \hookrightarrow for any node i° , we will be considering the i^{th} row of the adjacency matrix (A). Basically the feature of node i° will be the aggregation of the neighbourhood.

If $A_{ij} = 1$ then it belongs to the neighbourhood of i^{th} node

$$\text{if } j=1 \text{ to } N \quad A_{i^{\circ} j}$$

$$\text{Degree}(i^{\circ}) = \sum_{j=1 \text{ to } N} A_{i^{\circ} j}$$

- ④ If degree is very high too many things will be added
 - ⑤ If degree is low then very less things will be added
- } This may lead to vanish/explode the absolute value of the feature

As optimizers are quiet sensitive to I/P scales. \hookrightarrow This may result in backpropagation issues

\Downarrow Hence the finally aggregated node feature need to be suitably normalized.

Solution ②

The node features are normalized
by the degree of the node.
(Aggregated feature per unit)
of degree

$$D = \begin{bmatrix} d_1 & & & \\ & d_2 & & \\ & & d_3 & \dots \\ & & & d_N \end{bmatrix}$$

$$D^{-1} = \begin{bmatrix} \frac{1}{d_1} & & & \\ & \frac{1}{d_2} & & \\ & & \frac{1}{d_3} & \dots \\ & & & \frac{1}{d_N} \end{bmatrix}$$

Hence the final update rule can be

$$f(X, A) = D^{-1} A X \quad (a)$$

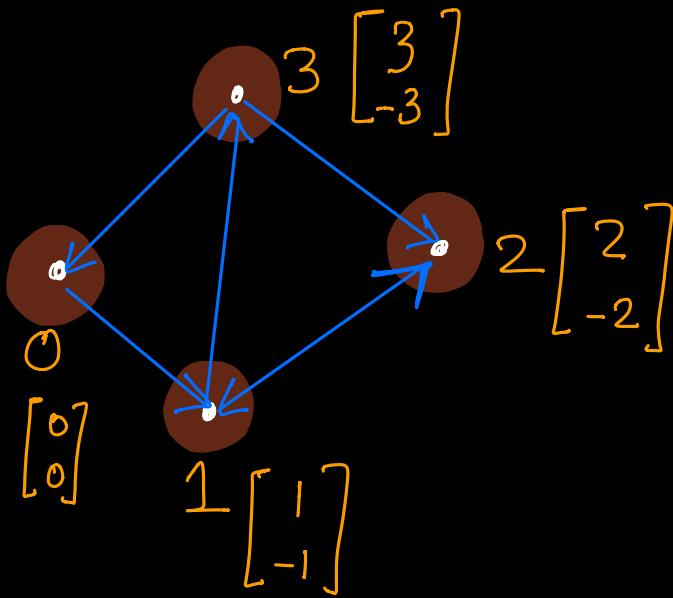
$$D^{-1} \hat{A} X \quad (b)$$

(Solving both issues)

$$\hat{A} = A + I$$

Therefore for our case

$$A = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \end{bmatrix}$$



$$X = \begin{bmatrix} 0 & 0 \\ 1 & -1 \\ 2 & -2 \\ 3 & -3 \end{bmatrix}$$

$$AX = \begin{bmatrix} 1 & -1 \\ 5 & -5 \\ 1 & -1 \\ 2 & -2 \end{bmatrix}$$

$$\hat{A} = A + I = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 0 & 1 & 1 & 1 \\ 0 & 1 & 1 & 0 \\ 1 & 0 & 1 & 1 \end{bmatrix}$$

$$\hat{A} X = \begin{bmatrix} 1 & -1 \\ 6 & -6 \\ 3 & -3 \\ 5 & -5 \end{bmatrix}$$

$$D = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 2 \end{bmatrix}$$

$$D^{-1} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0.5 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0.5 \end{bmatrix}$$

$$(a) D^{-1} \hat{A} X$$

$$\begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0.5 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0.5 \end{pmatrix} \begin{pmatrix} 1 & -1 \\ 5 & -5 \\ 1 & -1 \\ 2 & -2 \end{pmatrix} = \begin{pmatrix} 1 & -1 \\ 2.5 & -2.5 \\ 1 & -1 \\ 1 & -1 \end{pmatrix}$$

$$(b) D^{-1} \hat{A} X$$

$$\begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0.5 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0.5 \end{pmatrix} \begin{pmatrix} 1 & -1 \\ 6 & -6 \\ 3 & -3 \\ 5 & -5 \end{pmatrix} = \begin{pmatrix} 1 & -1 \\ 3 & -3 \\ 3 & -3 \\ 2.5 & -2.5 \end{pmatrix}$$

This aggregation utilizes
nodes' feature as well as aggregated.

✳ Now learning the transformation parameterized by (W)

$$f_W = (D^{-1} A^T X) * W$$

$$(a) f_{W_1} = (D^{-1} A^T X) * W_1$$

$$\begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}$$

$$\begin{bmatrix} 1 & -1 \\ 3 & -3 \\ 3 & -3 \\ 2.5 & -2.5 \end{bmatrix}_{4 \times 2} \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}_{2 \times 2} = \begin{bmatrix} 2 & -2 \\ 6 & -6 \\ 6 & -6 \\ 5 & -5 \end{bmatrix}_{(4 \times 2)}$$

$$(b) f_{W_2} = (D^{-1} A^T X) * W_2$$

$$\begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

In order to reduce the feature's dimensionality

$$\begin{bmatrix} 1 & -1 \\ 3 & -3 \\ 3 & -3 \\ 2.5 & -2.5 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \end{bmatrix} = \begin{bmatrix} 2 \\ 6 \\ 6 \\ 5 \end{bmatrix}_{(4 \times 1)}$$

✳ After aggregation & transformation, nonlinearity has been introduced using ReLU Activation (σ)

$$\text{ReLU}(f_W(X, A))$$

(a)

$$\begin{bmatrix} 2 & 0 \\ 6 & 0 \\ 6 & 0 \\ 5 & 0 \end{bmatrix}$$

These are the complete final representation.

$$\begin{bmatrix} 2 \\ 6 \\ 6 \\ 5 \end{bmatrix}$$

$$H = \sigma(D^{-1} A^T X * W)$$

Previous layer O/P.

∴ for Aggregation ($A, X, \& D$) are required
Transformation ($W \& \sigma$) are required

Hence one can see that any function (f_w) learned at any layer (l), can be represented as

$$f_{W^t}(H^t, A) = \text{Transform}_{//}(\text{Aggregate}_{//}(A, H^t), W^t)$$

Transform(M, W)
= $\mathcal{O}(M * W)$

Rule (1) (sum)

Aggregati (A, H^l) = AH

Rule (2) (MEAN)

$$\text{Aggregat } (A, H) = D^{-1} \hat{A}^H \\ (\hat{A}^H = A + I)$$

Basically both aggregation rule uses weighted sum of its neighbouring features. Just weighing is different.

A) SUM Rule (Neighbourhood Aggregation)

$$\begin{array}{c}
 \text{3x3} \\
 A_1 \begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{pmatrix} \times \begin{pmatrix} a & b & c \\ d & e & f \\ g & h & i \end{pmatrix} = \begin{pmatrix} 4a + 5d + 6g \\ 4b + 5e + 6h \\ 4c + 5f + 6i \end{pmatrix}^T \\
 \text{3x3} \\
 A_2 \begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{pmatrix} \times \begin{pmatrix} a & b & c \\ d & e & f \\ g & h & i \end{pmatrix} = \begin{pmatrix} 4a + 5d + 6g \\ 4b + 5e + 6h \\ 4c + 5f + 6i \end{pmatrix}^T \\
 \text{3x3} \\
 A_3 \begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{pmatrix} \times \begin{pmatrix} a & b & c \\ d & e & f \\ g & h & i \end{pmatrix} = \begin{pmatrix} 4a + 5d + 6g \\ 4b + 5e + 6h \\ 4c + 5f + 6i \end{pmatrix}^T
 \end{array}$$

Let us working on i^{th} node, hence using i^{th} row of adjacency matrix (A_i) .

$$\text{agg}(A, X)_i = A_i^T X = \sum_{j=1}^N A_{ij} * X_j$$

↓
 only included if it's a neighbour
 feature of (j^{th}) node

Hence Contribution of each Node depends solely on neighbourhood defined by adjacency matrix (A) .

B) Mean Rule

$$\begin{aligned} \text{agg}(A, X)_i &= D^{-1} A_i^T X \\ &= \sum_{k=1}^N D_{i,k}^{-1} * \sum_{j=1}^N A_{ij} * X_j \\ &\quad \text{Just same as previous.} \\ &= \sum_{j=1}^N D_{i,j}^{-1} * A_{ij} * X_j = \sum_{j=1}^N \frac{A_{ij}}{D_{ii}} * X_j \\ &\quad \text{Contribution got weighted by degree of the } (i^{th}) \text{ node} \end{aligned}$$

(i) is fixed &
 (k) varies from 1 to N

feature of j^{th} node

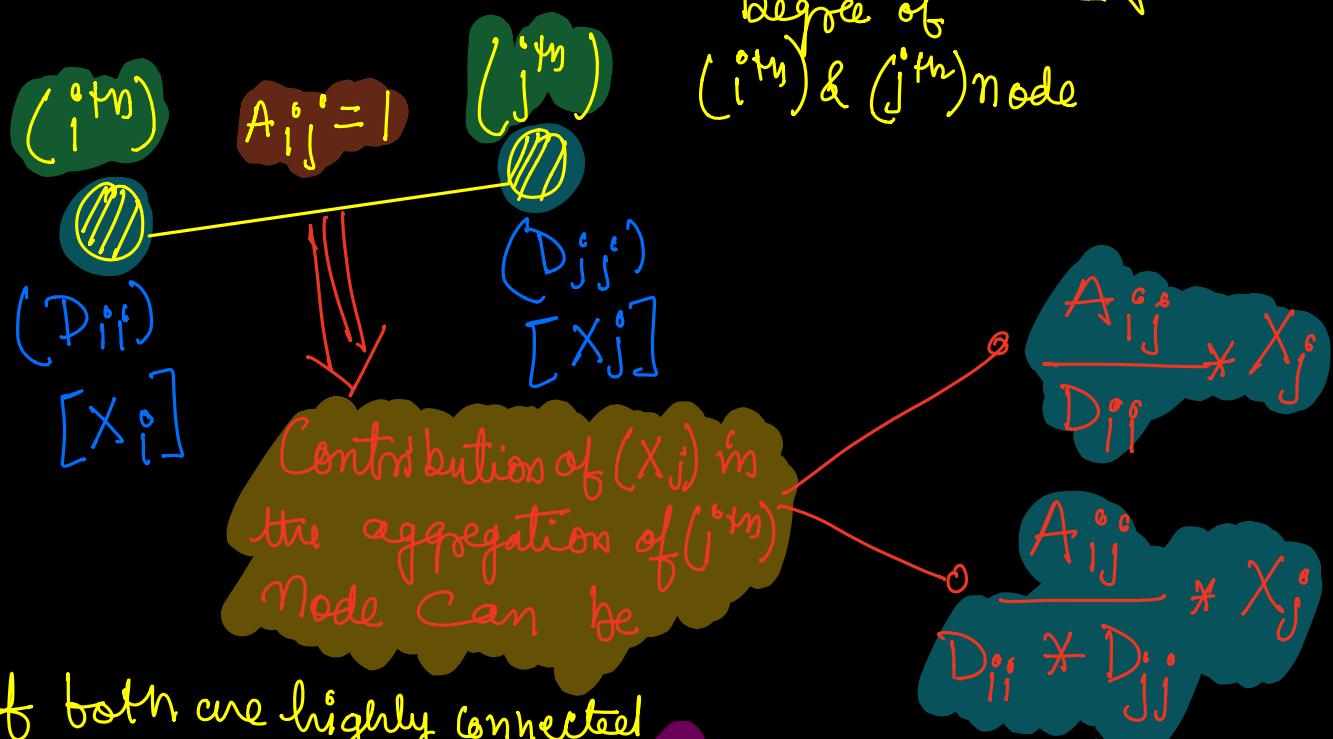
$$\sum_{j=1}^N \frac{A_{ij}}{D_{ii}} = \frac{\sum_{j=1}^N A_{ij}}{D_{ii}} = (1)$$

Weights are guaranteed to be summed equal to (1).
for any node (i) [Weighted Normalization]

This can be seen as the mean feature representation of neighbours of i^{th} mode.

C) Spectral Rule

$$\text{agg}(A, X)_i = D^{-\frac{1}{2}} A_i D^{-\frac{1}{2}} X = \sum_{j=1}^N \frac{A_{ij}}{\sqrt{D_{ii} D_{jj}}} X_j$$



* If both are highly connected
Less Contribution ↓

* If both are Sparsely Connected
High Contribution ↑