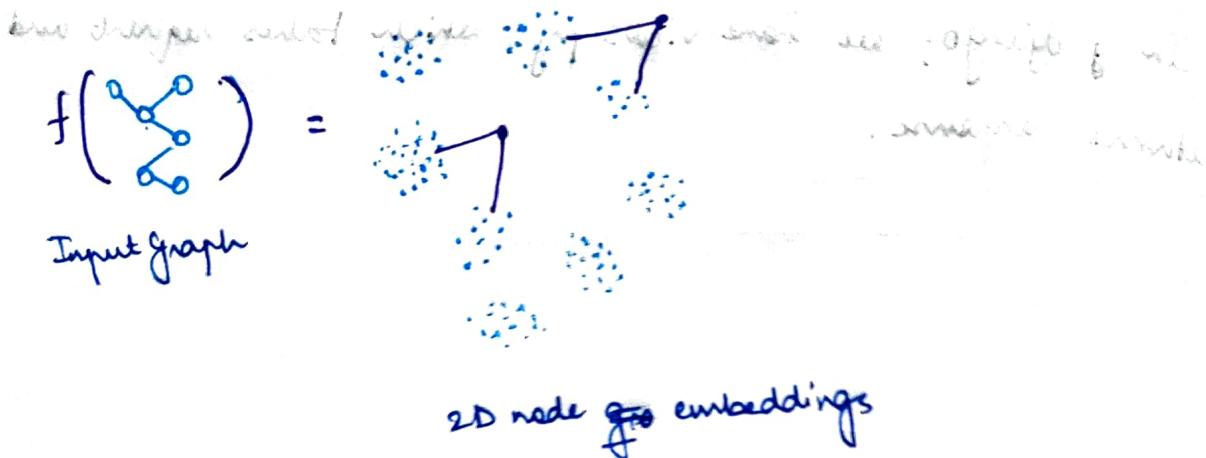


Node Embeddings

Intuition: Map nodes to d -dimensional embeddings, such that similar nodes in the graph are embedded close together.



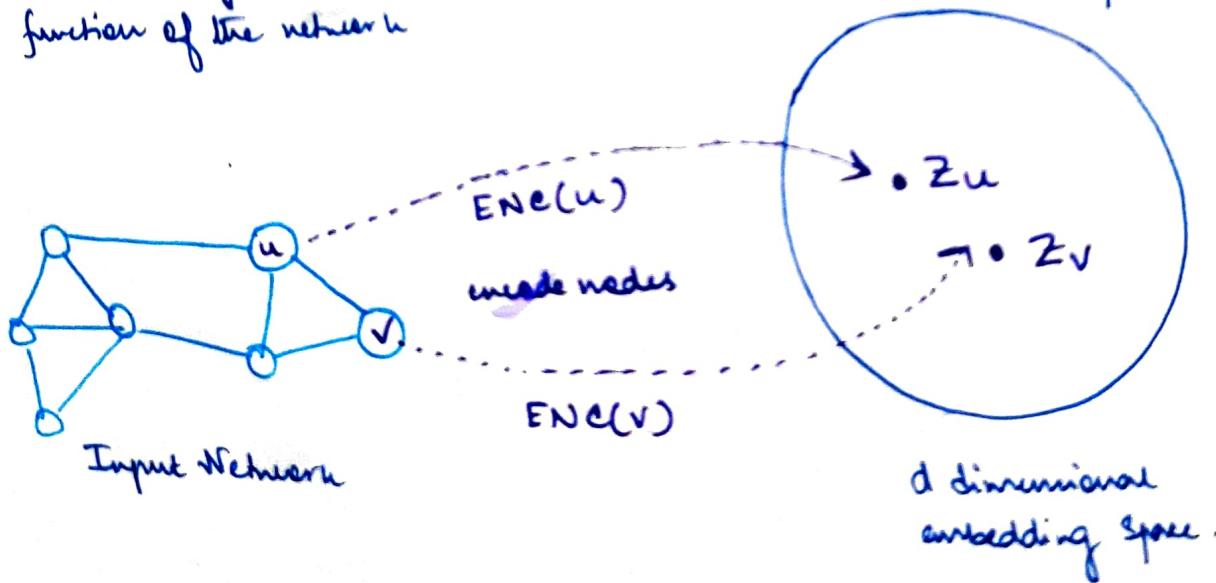
The idea is to encode nodes in such a way that

$$\text{similarity}(u, v) \approx z_v^T z_u$$

$$\boxed{\text{similarity}(u, v) \approx z_v^T z_u}$$

similarity function of the network

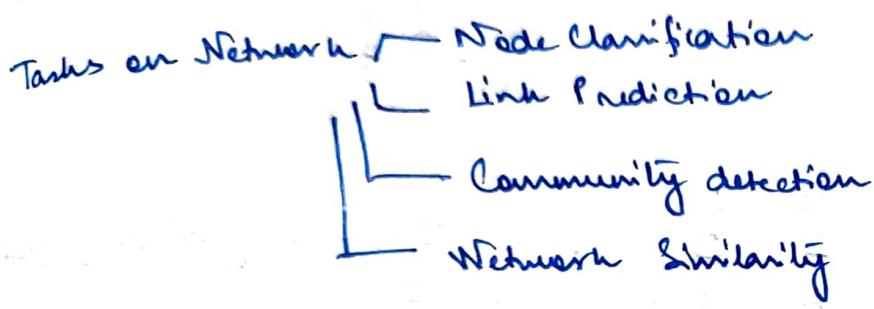
similarity function of the vector space



so, if two nodes are connected by edges, then they should be close together in the embedding space.

In GNN, for a node v

$\text{ENC}(v) = \text{multiple layers of}$
 $\text{non-linear transformations}$
 $\text{based on graph structure}$



Deep learning for Graphs

From what we have in classical deep learning, the question comes
HOW TO GENERALIZE CLASSICAL DEEP LEARNING TO WORK WITH
GRAPHS?

Set Up: Assume we have a graph $G(V, E)$.

V : Vertex Set

$E: E \subseteq V \times V$ is the set of edges

$A: A_{|V| \times |V|}$ is the adjacency matrix

$X: X \in \mathbb{R}^{|V| \times d}$ is the matrix of node feature.

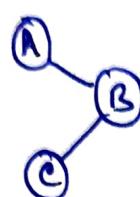
v : a given node of G .

N_v : (or) $N(v)$ is the ~~set of~~ neighbourhood set of the node v .

The most important question is:

HOW DO WE INPUT A GRAPH INTO A NEURAL NETWORK?

The Naive Approach is to use the adjacency matrix of the graph, concatenate it with corresponding node features and feed it to the graph network.

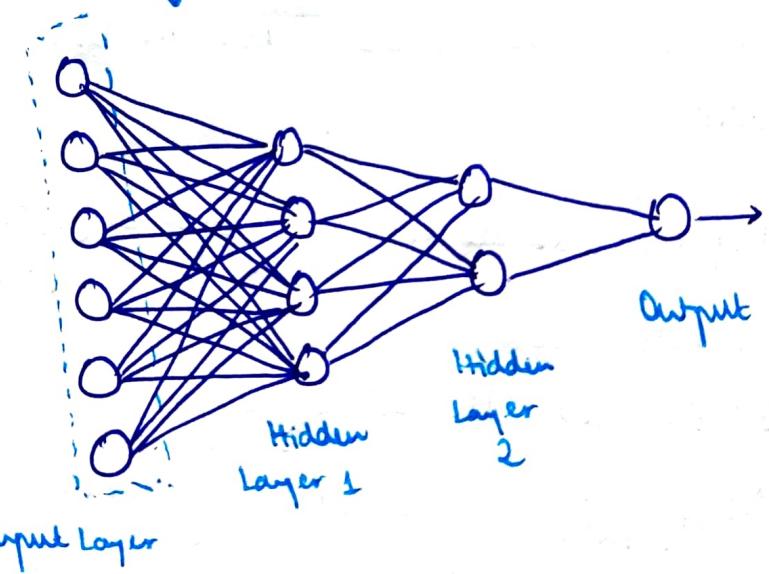


$$A \rightarrow \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

$$B \rightarrow \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

$$C \rightarrow \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

$$\begin{array}{c} A \quad B \quad C \\ \hline A \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} \quad \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \rightarrow [\text{Adj}(A), \text{feat}(A)] \end{array}$$



Issues:

- Since, the input is a combination of the adjacency entry of a node and its features the number of parameters is $O(|V|)$ and since, we have one training example per node

The number of parameters would be greater than the training examples, resulting in unstable training and overfitting.

- This way of ~~encoding~~ inputting a graph in a neural network will not be applicable for graphs with different sizes since, that would have different adjacency matrices.
- This method is sensitive to node ordering, as the adjacency matrix would change if we label the nodes in a different way. This is result in a different input being passed to the neural network and thereby confusing the model.

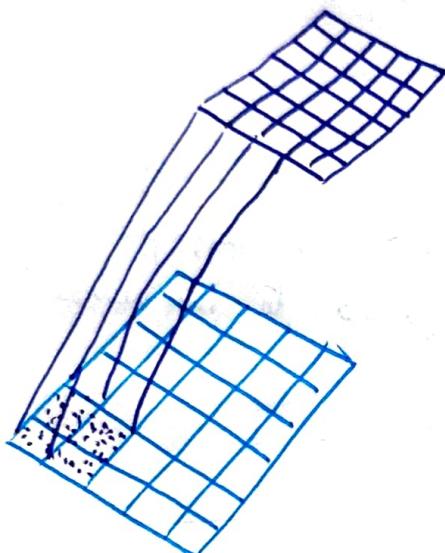
Idea: Convolutional Networks

The idea is to generalise convolutions beyond simple lattices and use it on graphs leveraging the node features / attributes.

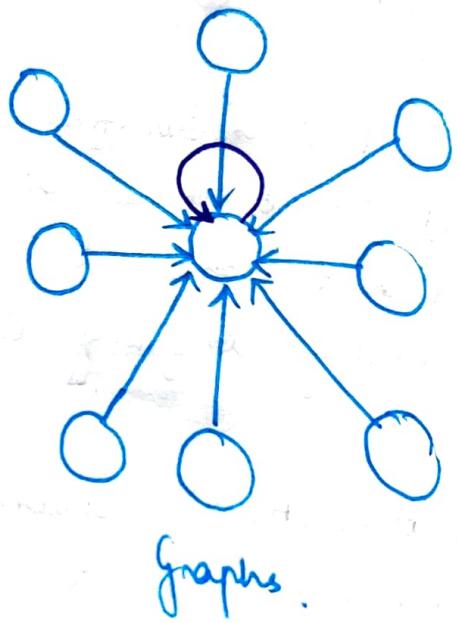
Problems with using the same convolutional operator (as used in images) in graphs:

- No fixed notion ~~vector~~ of location or sliding window on the graph.
- Graph is permutation invariant i.e. the properties of graphs don't change by permuting the node ordering. So, we need a function / model that doesn't depend on the arbitrary ordering of the adjacency matrix.

From Images to Graphs



Image



Graphs.

The idea of a convolution is, it takes a submatrix of pixel applies some transformation on it to generate a new pixel and then we slide this operator over the entire image.

In graphs, to incorporate this technique we have a ^{central} node which receives information from its neighbours and creates new information using them as well as its own information.

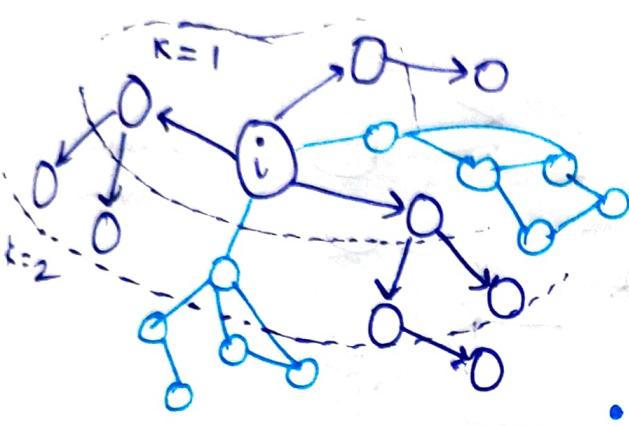
Idea: Transform information from neighbours and combine it:

- Transform messages h_i from neighbours: $W_i h_i$
- Add them up $\sum_i W_i h_i$

~~How to feed a G~~

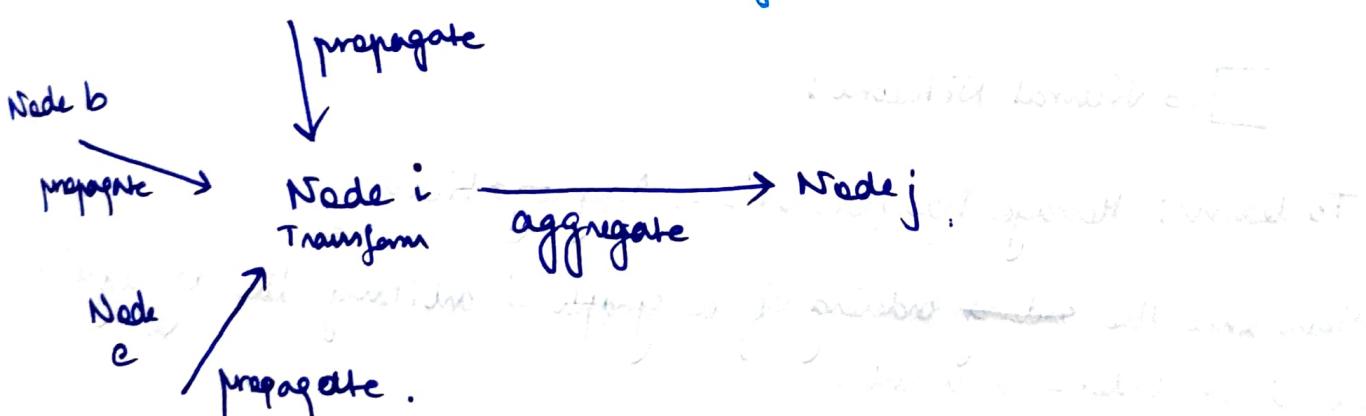
How graph convolutional Networks work?

In GNN's every node has its own neural network defined by its neighbourhood.



So, in order to just predict something about node i, then the node i will take information from its neighbours, which will take information from its neighbours and learn to

- Propagate this information across the graph
- Aggregate this information
- Create a new message which the next node can propagate, transform and aggregate.

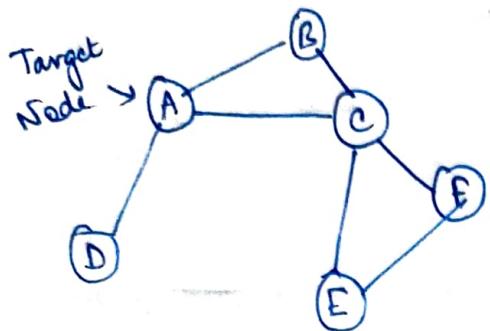


So, the idea in a GNN is just that kind of a computation graph.

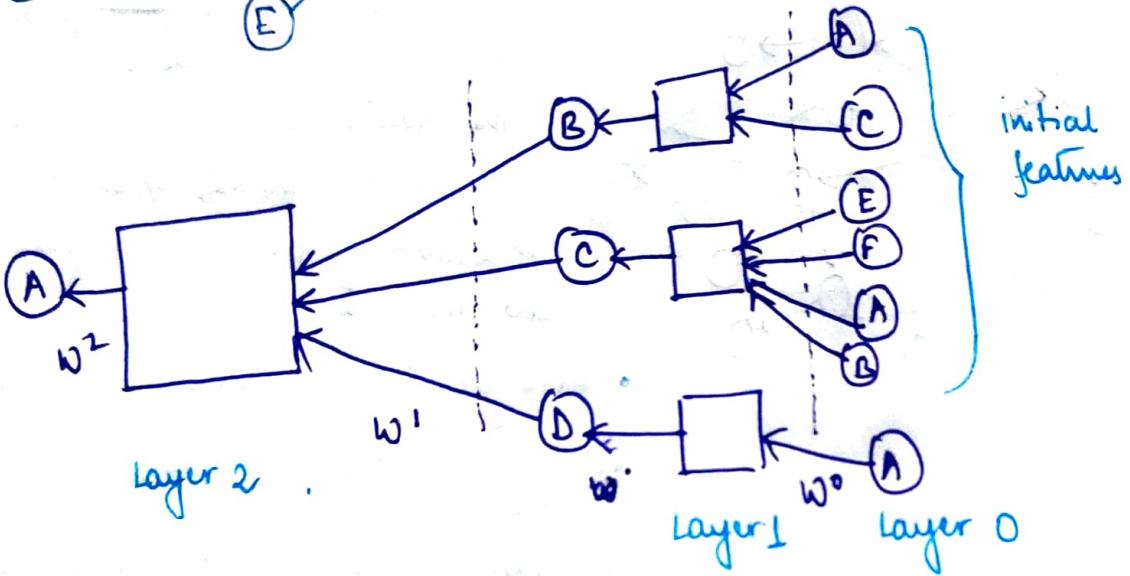
- Determine the node computation graph (local neighbourhood)
- Propagate and transform this information over the computation graph.

This computation graph determines the underlying structure of the neural network.

suppose we have a graph G as follows:



So, A will have its own neural network architecture (clearly) which would look as follows



Network Architecture for layer A .

\square := Neural Network :

To learn: Message Transformation, Aggregation.

Now, since the ~~existing~~ ordering of a graph is arbitrary the aggregation should be order-invariant.

NEIGHBORHOOD AGGREGATION SHOULD BE
PERMUTATION INVARIANT

• simply neighborhood aggregation would be

- Average the information from messages from the neighbors

• Apply the neural network.

mathematically,

$h_v^l :=$ embedding of node v in layer l

$W_L :=$ weight matrix in layer L for neighbours

$B_L :=$ weight matrix in layer L for the original node (v).

$$h_v^0 = x_v \rightarrow \text{node feature of node } v$$

$$h_v^{(L+1)} = \sigma \left(\frac{1}{|N(v)|} \sum_{u \in N(v)} [W_L \cdot h_u^{(L)} + B_L h_v^{(L)}] \right)$$

$z_v = h_v^{(L)}$ \rightarrow final node embedding $\in \{0, 1, \dots, L-1\}$

$\sigma :=$ transformation function (non linear e.g ReLU).

Now, if every node has its own architecture and transformation operator, the question comes

HOW DO WE TRAIN THE MODEL?

The parameters are W_L and B_L $\in \{0, 1, \dots, L-1\}$ are trainable and can be trained using SGD by feeding their embeddings in any loss function.

i.e let $\hat{y} := f(h_v^L)$

$y :=$ target label for node v .

Then,

$$W^* = \underset{W^0, W^1, \dots, W^{L-1}}{\operatorname{argmin}} \ell(\hat{y}, y)$$

can be solved using SGD.

Host Matrix formulation:

Let $H^{(L)}$

$$\text{Let } H^{(L)} = [h_1^{(L)} \ h_2^{(L)} \ \dots \ h_{|V|}^{(L)}]^T$$

Now, $\sum_{u \in N(v)} h_u^{(L)}$ is basically $A_v H^{(L)}$

Let $D \in \mathbb{R}^{N \times N}$ be a diagonal matrix s.t

$$D_{v,v} = \deg(v)$$

$$D_{u,v} = 0 \text{ if } u \neq v.$$

Then,

$$DD^{-1} = I$$

$\rightarrow D^{-1} = D^{-1}$ is also a diagonal matrix with entries:

$$D_{v,v}^{-1} = \frac{1}{\deg(v)}$$

$$D_{u,v}^{-1} = 0 \text{ if } u \neq v$$

$$\text{So, } \frac{1}{|\mathcal{N}(v)|} \sum_{u \in N(v)} h_u^{(L)} \text{ is } \cancel{\sum_{u \in N(v)}} \sum_{u \in N(v)} \frac{h_u^{(L)}}{|\mathcal{N}(v)|}$$

which is:

$$D_{v,v}^{-1} \sum_{u \in N(v)}$$

$$D_{v,v}^{-1} A_v H^{(L)}$$

∴ generalizing this

$$\frac{1}{|N(v)|} \sum_{u \in N(v)} h_u^{(l)} = D^{-1} A H^{(l)}$$

Rewriting the update function in matrix form:

$$H^{(l+1)} = \sigma \left(D^{-1} A H^{(l)} W_l^T + H^{(l)} B_l^T \right)$$

$$H^{(l+1)} = \sigma \left(\tilde{A} H^{(l)} W_l^T + H^{(l)} B_l^T \right) \quad \text{where } \tilde{A} = D^{-1} A$$

neighborhood
aggregation

self
transformation

How to train a GNN?

Supervised setting : Suppose we are interested in node classification and so, in this supervised setting we minimize the loss

$$\min_{\theta} \ell(y, f(z_v))$$

y := node label of v

z_v := embedding of node v

$$\min_{\theta} \ell(y, f(z_v))$$

We can use SGD to come to the solution.

~~step~~ Unsupervised setting: In the unsupervised setting the graph structure can be used as supervision.

We can define the notion of similarity as the dot product of two nodes should between the two node embeddings should correspond to their similarity in the network.

Now, we can use Deep Encoder as the to learn the embeddings.

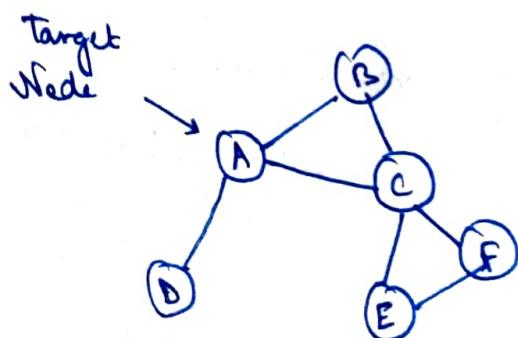
$$L = \sum_{z_u, z_v} CE(y_{u,v}, DEC(z_u, z_v))$$

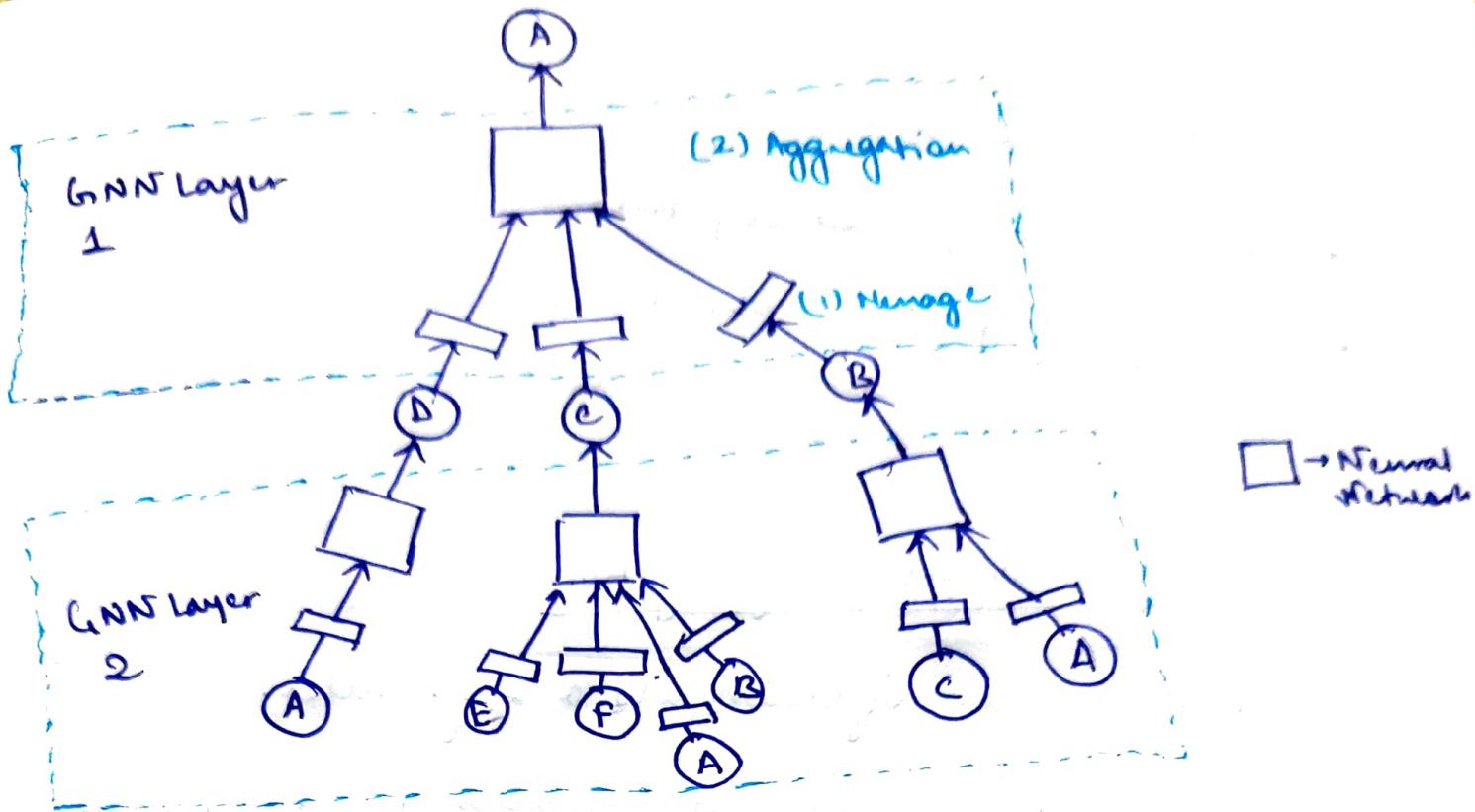
$y_{u,v} = 1$ if u and v are similar.

$DEC(u, v)$ = decoder such as the inner product or cosine similarity.

A General GNN Framework

A General GNN Framework has three components to talk about.
Consider we have a graph





The first component of a GNN framework is:

- Merge + Aggregation which is the content of every layer of a GNN

The second component is about:

- Stacking layers up in a GNN (correct GNN layers)
 - We can stack layers sequentially
 - We can add skip connections

third

The last component is:

- Defining the computational graph. Ideas like
Raw Input Graph ≠ Computational graph

Ideas like:

- Graph feature augmentation
- Graph structure manipulation

The final component is:

How to train the GNN?

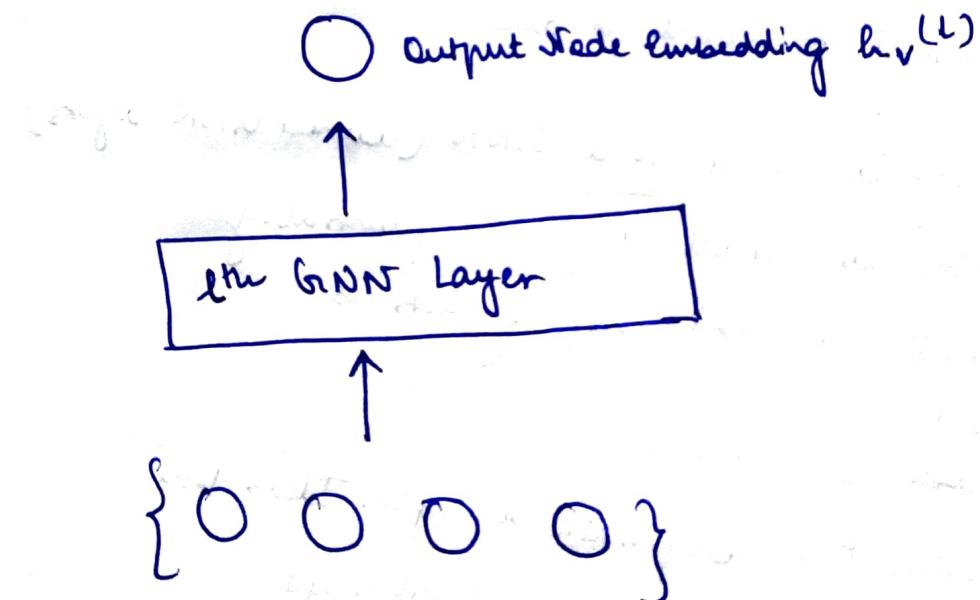
- Supervised / Unsupervised objectives
- Node / Edge / Graph level objective.

A GNN Layer : Message + Aggregation

Idea: compress a set of vectors into a single vector.

We want to compress messages coming from the neighbors
and we do this by a two step process:

- Message Transformation
- Message Aggregation



Input: $h_v^{(L-1)}$ and $h_u^{(L-1)}$; $u \in N(v)$.

- **Message Computation:** In this step, we take the output from the previous layer and somehow transform it to create a message to be sent to other nodes layer.

$$m_u^{(L)} = \text{MSCh}(h_u^{(L-1)})$$

Typically this message computation is a linear layer described as follows

$$m_u^{(L)} = W_1 h_u^{(L-1)}$$

- **Aggregation:** This is the step, where the transformed messages are aggregated into a single message.

$$h_v^{(L)} = \text{AGCh}(\{m_u^{(L)} \mid u \in N(v)\})$$

The AGCh function can be the Sum(\cdot) , Max(\cdot) or Mean(\cdot) aggregator.

Issue: The issue with is that the information from the even node v could get lost and as we are just using the messages from the neighbors.

Solution: Include $h_v^{(L-1)}$ while computing $h_v^{(L)}$

- **Message:**

$$\cancel{m_u} = \text{MSCh } m_u = W_1 h_u^{(L-1)}$$

message from neighbors

$$m_v = \text{BCh } h_v^{(L-1)}$$

message from even node

- Aggregation: After aggregating the messages from neighbors we aggregate the message from itself via summation or concatenation.

$$w_v(l) = \text{concat}\left(\text{AGG}\left(\left\{h_u^{(l-1)} | u \in N(v)\right\}\right)\right),$$

So, putting things together,

(1) Manage: Each node computes a manage,

$$m_u^{(l)} = \text{NSu}(m_u^{(l-1)}) \quad \text{if } u \in N(v) \cup \{v\}$$

(2) Aggregation: Aggregate messages from neighbors

$$z_v^{(l)} = \text{arg} \left(\left\{ m_u^{(l)} \mid u \in N(v) \right\}, m_v^{(l)} \right)$$

(3) Transformation: Finally, we can transform the ~~other~~ output to add expressiveness to our model.

$$h_v(l) = \sigma(z_v(l))$$

σ : ReLU, Sigmoid, etc.

Canonical GNN Layers : GCN (Graph Convolutional Networks)

for GCN:

$$h_v^{(L)} = \sigma \left(W^{(L)} \sum_{u \in N(v)} \frac{h_u^{(L-1)}}{|N(v)|} \right)$$

writing this in the message + aggregation form:

$$h_v^{(L)} = \sigma \left(\sum_{u \in N(v)} W^{(L)} \frac{h_u^{(L-1)}}{|N(v)|} \right)$$

Graph SAGE :

$$h_v^{(L)} = \sigma \left(W^L \cdot \text{concat} \left(h_v^{(L-1)}, \text{Agg} \left(\{ h_u^{(L-1)} \mid u \in N(v) \} \right) \right) \right)$$

Writing it as Message + Aggregation form:

- Message is computed ~~through~~ within $\text{Agg}(\cdot)$

Two Stage Aggregation:

Stage 1: Aggregate from node neighbors

$$h_{N(v)}^L = \text{Agg} \left(\{ h_u^{(L-1)} \mid u \in N(v) \} \right)$$

Stage 2: Further aggregate over the node itself

$$h_v(l) = \sigma(w(l), \text{concat}(h_v^{(l-1)}, h_N^l(v)))$$

Aggregation Functions:

Mean: Take weighted average of neighbors

$$\text{AAG} = \sum_{u \in N(v)} \frac{h_u^{(l-1)}}{|N(v)|}$$

Pool: We can transform the neighbor vectors and apply Mean(.) or Max(.) .

$$\text{AAG} = \text{Mean}(\{\text{MLP}(h_u^{(l-1)}) \mid u \in N(v)\})$$

Linear or non linear transformation.

LSTM: We can even apply LSTM to the neighbors. However, since, LSTM's are not order invariant, we must feed it a permutation of the ~~most~~ neighbors everytime so as to teach it not to take the order into consideration.

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

The choice of the aggregation function ~~decides~~ decides how expressive our GNN would be.

$$h_v^{(l)} = \sigma(w^{(l)}. \text{concat}(h_v^{(l-1)}, h_{N(v)}^l))$$

Aggregation Functions:

Mean: Take weighted average of neighbors

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

Pool: We can transform the neighbor vector and apply Mean(.) or Max(.) .

$$ACh = \text{Mean}(\{\text{MLP}(h_u^{(l-1)}) \mid u \in N(v)\})$$

Linear or non linear transformation.

LSTM: We can even apply LSTM to the neighbors. However, since, LSTM's are not order invariant, we must feed it a permutation of the ~~most~~ neighbors everytime so as to teach it not to take the order into consideration.

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

The choice of the aggregation function ~~decides~~ decides how expressive our GNN would be.

L_2 Normalization:

In GraphSAGE, L_2 Normalization is applied at every layer. That is,

$$h_v^{(L)} = \frac{h_v^{(L)}}{\|h_v^{(L)}\|_2}$$

$\forall v \in V$

- Without L_2 normalization, the embedding vectors have different scales.
- In some cases, normalization also result in performance improvement.

Graph Attention Networks

$$h_v^{(L)} = \sigma \left(\sum_{u \in N(v)} w^{(L)} \alpha_{uv} h_u^{(L-1)} + \alpha_{vv} h_v^{(L-1)} \right)$$

$$h_v^{(L)} = \sigma \left(\sum_{u \in N(v)} w^{(L)} \alpha_{uv} h_u^{(L-1)} \right)$$

↑
Attention weights.

In comparison to GraphSAGE and GCN these attention weights α_{uv} were basically:

$$\alpha_{uv} = \frac{1}{|N(v)|} \quad \text{and these were based on the structural properties of graphs like node degree.}$$

Also clearly in this all neighbors $u \in N(v)$ are equally important

to node v.

Now, the idea in GATs is

Not all node's neighbors are equally important.

- Attention inspired by cognitive attention focuses on the important parts of the input data and throws out the rest.
- Which part of the data is more important depends on the context and is learned through training.

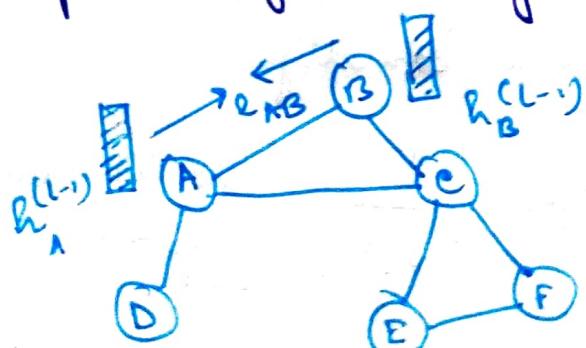
So, can we let the weighing factors α_{uv} to be learned?

Then, weighing factors are learned along with all the other parameters of the network.

Then, coefficients α_{uv} are learnt as a byproduct of an attention mechanism ~~as~~ α .

$$\alpha_{uv} = \text{a}\left(W^{(l)}h_u^{(l-1)}, W^{(l)}h_v^{(l-1)}\right)$$

α_{uv} indicates the importance of u's message to node v.



So, some function α is defined that will take the transformed embedding of u at previous layer, embedding of v at previous layer transform these into messages and then return an attention score. $\alpha \in \mathbb{R}$

Then, we normalize these weights by α_{uv} using the softmax function.

$$\alpha_{uv} = \frac{\exp(\alpha_{uv})}{\sum_{w \in N(v)} \exp(\alpha_{uw})}$$

$$\boxed{\alpha_{uv} = \frac{\exp(\alpha_{uv})}{\sum_{k \in N(v)} \exp(\alpha_{kv})}}$$

Now, the form of attention mechanism is dependent on the choice of α . A simple choice of α would be to concatenate the two followed by a linear layer where the weights of the linear layer would be jointly learnt with the weights of the RNN.

$$\alpha_{uv} = \text{Leaky ReLU}\left(a^T \cdot [w^{(l)}_{hu} h_u^{(l-1)}, \| w^{(l)}_{hv} h_v^{(l-1)}] \right)$$

$$a^T \in \mathbb{R}^{2F} \quad \text{where, } w^{(l)}_{hu} h_u^{(l-1)} \in \mathbb{R}^F$$

$\|$: concatenation operation

$$w^{(l)}_{hv} h_v^{(l-1)} \in \mathbb{R}^F$$

$$\alpha_{uv} = \frac{\exp(\text{LeakyReLU}(a^T [w^{(l)} h_u^{(l-1)} || w^{(l)} h_v^{(l-1)}]))}{\sum_{k \in N(v)} \exp(\text{LeakyReLU}(a^T [w^{(l)} h_k^{(l-1)} || w^{(l)} h_v^{(l-1)}]))}$$

Now, working with these attention weights can be tricky as they can be hard to learn or they may be hard to converge.

So, in order to stabilize learning we introduce Multihead attention which stabilizes the learning process.

In order to do that we employ obtain multiple attention scores and use it to compute multiple embeddings of the same node. Then, we concatenate them.

$$h_v^{(l)}[1] = \sigma \left(\sum_{u \in N(v)} w^{(l)} \alpha_{uv}^1 h_u^{(l-1)} \right)$$

$$h_v^{(l)}[2] = \sigma \left(\sum_{u \in N(v)} w^{(l)} \alpha_{uv}^2 h_u^{(l-1)} \right)$$

$$h_u^{(l-1)} \in \mathbb{R}^F$$

$$h_v^{(l)}[k] = \sigma \left(\sum_{u \in N(v)} w^{(l)} \alpha_{uv}^k h_u^{(l-1)} \right)$$

$$h_v^{(L)} = \underset{i=1}{\overset{K}{||}} h_v^{(L)}[i]$$

$||$ represents concatenation.

So, instead of F features we now have KF features of these embeddings.

Benefits of Attention Mechanism:

- Key Benefit: Allows specifying different importance values to different neighbors
- Computationally Efficient: Computation of attention coefficients can be parallelized across all edges.

Aggregation may be parallelized across all nodes.

Storage efficient:

- Sparse Matrix operations don't require more than $O(V+E)$ entries to be stored.
- Fixed number of parameters irrespective of graph size (as, the weight in the layer of attention mechanism depend on the features of node embeddings and not the graph).

Localized:

Only attends to local network neighborhoods.

Inductive capability:

It is shared edge-wise mechanism
Doesn't depend on the global graph structure.

Modifying all of these GNN architectures to include their own self representations of nodes while learning embeddings.

GCN:

$$h_v^{(0)} = x_v \quad \forall v \in V$$

$$h_v^{(l)} = \sigma \left(W^l \cdot \frac{1}{|N(v)|} \sum_{u \in N(v)} h_u^{(l-1)} + B^l h_v^{(l-1)} \right) \quad \forall v \in V$$

GraphSAGE:

~~$$h_v^{(l)} = \sigma \left(W^l \cdot \text{Agg}_{u \in N(v)} \left(\{h_u^{(l-1)}\} \right), h_v^{(l-1)} \right)$$~~

$$h_v^{(0)} = x_v \quad \forall v \in V$$

$$h_v^{(l)} = \sigma \left(W^l \cdot \left[\text{Agg}_{u \in N(v)} \left(\{h_u^{(l-1)}\} \right), h_v^{(l-1)} \right] \right) \quad \forall v \in V.$$

GAT

$$h_v^{(0)} = x_v \quad \forall v \in V.$$

$$h_v^{(l)} = \sigma \left(W^{(l)} \cdot \left[\sum_{u \in N(v)} \alpha_u^{(l)} h_u^{(l-1)} + \alpha_v^{(l)} h_v^{(l-1)} \right] \right)$$