



## Personal Note #2: Graph Neural Network.

→ Graph Neural Networks work with Graph-like data.

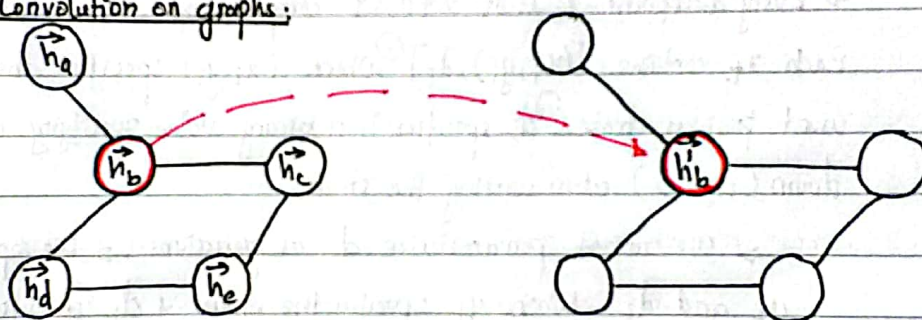
### ② Graph convolution.

→ Graph can be seen as strict generalisation of images.

→ CNN leverages convolutional operations to extract the spatial regularity in images.

→ Would be nice if we could generalize it to operate on arbitrary graphs!

### \* Convolution on graphs:



$$\vec{h}'_i = g(\vec{h}_a, \vec{h}_b, \dots, \vec{h}_e).$$

→ Ideally  $g(\cdot)$  should be some function of the entire set of neighbors to particular nodes.

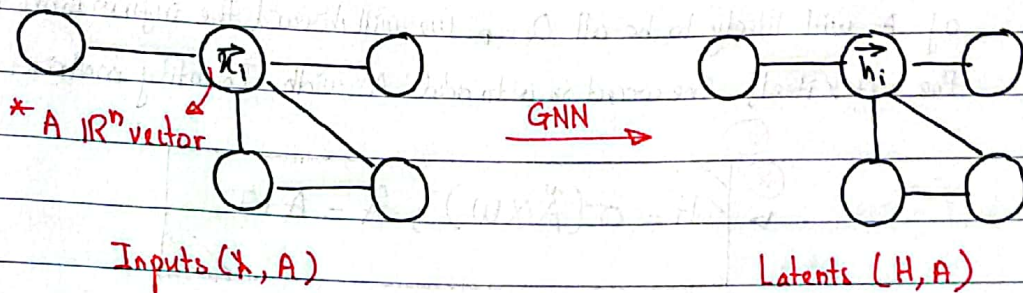
→ ~~Challenges with graph convolutions.~~ Desired properties for graph conv:

- Computational efficiency ( $\sim O(V+E)$ )
- Fixed number of params (independent of inputs).
- Localisation (act on local neighborhood of nodes).
- Weighted neighbors importance.
- Inductivity.

→ However, arbitrary graphs are much harder to deal with.



\* What to do with GNN outputs.



- \* {
- ① Node classification:  $\vec{z}_i = f(\vec{h}_i)$  \*
  - ② Graph classification:  $\vec{z}_G = f(\sum \vec{h}_i)$
  - ③ Link classification:  $\vec{z}_{ij} = f(\vec{h}_i, \vec{h}_j, \vec{e}_{ij})$

\* Simple graph update rule:

↳ Let's assume our graph is unweighted and undirected:

$$A_{ij} = A_{ji} = \begin{cases} 1, & \text{if } x_i \text{ and } x_j \text{ are connected.} \\ 0, & \text{otherwise.} \end{cases}$$

↳ We can aggregate the neighborhood by multiplying with adjacency matrix:

①  $H = \sigma(AXW)$

→ {

- $X$ : Input nodes.
- $A$ : Adj matrix
- $W$ : Learnable weights.

- $X \in \mathbb{R}^{N \times d}$  where {
  - $N = \#$  of nodes in the graph.
  - $d = \#$  of dimensions in the feature vector/node.
- $A \in \mathbb{R}^{N \times N}$  where nodes order in  $A$  is consistent with nodes order in  $X$ .
- $W \in \mathbb{R}^{d \times d'}$  where  $d' = \text{latent's dimensions}$ .



(\*)  $\rightarrow$  However, there is one issue in (1), Unless enforced, the diagonal of  $A$  will likely to be all 0s  $\rightarrow$  we will discard the information on the vertex itself. One correction is to add  $A$  with Identity matrix.

(2)

$$H = \sigma(\tilde{A}XW); \tilde{A} = A + I$$

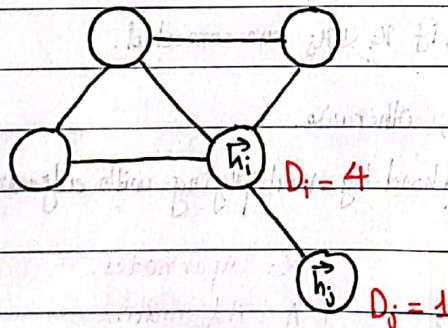
$\rightarrow$  The rule can be re-written node-wise as:

(2)

$$\vec{h}_i = \sigma \left[ \sum_{j \in N_i} W \vec{h}_j \right]; N_i = i^{\text{th}} \text{ node neighborhood}$$

(\*)

$\rightarrow$  Another problem is that the neighborhood size for every node is inconsistent  $\rightarrow$  The output scale might differ a lot.



$\rightarrow$  A simple remedy is to take the mean output at each node by dividing the output by the neighborhood size.

$\rightarrow$  We can normalize the adj matrix with the deg matrix  $D$ .

(4)

$$H = \sigma(\tilde{D}^{-1} \tilde{A}XW); \tilde{D} = \text{deg mat of } \tilde{A}.$$

(5)

$$\vec{h}_i = \sigma \left[ \frac{1}{|N_i|} \sum_{j \in N_i} W \vec{h}_j \right]$$



\* Other aggregation rules:

↳ Symmetric pooling:

$$H = \sigma(\tilde{D}^{-\frac{1}{2}} A \tilde{D}^{-\frac{1}{2}} X W)$$

→ Node-wise:

$$\vec{h}_i = \sigma \left[ \sum_{j \in N_i} \frac{1}{\sqrt{|N_i| |N_j|}} W \vec{h}_j \right]$$

→ Simple and powerful → Most cited GNN paper!

\* Message Passing Neural Network

↳ Let  $\vec{m}_{ij}$  be the message sent across  $i \rightarrow j$ , computed using a message function:

$$\vec{m}_{ij} = f_e(\vec{h}_i, \vec{h}_j, \vec{e}_{ij})$$

\* Edge features

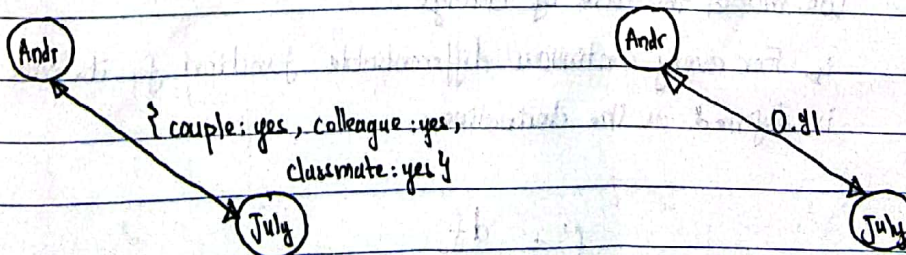
↳ A node's update rule is a function of its previous state and the aggregation of all the messages. This is called a readout function.

$$\vec{h}_i = f_v(\vec{h}_i, \sum_{j \in N_i} \vec{m}_{ij})$$

↳ Both  $f_e(\cdot)$  and  $f_v(\cdot)$  are small MLPs.

\* Examples of edge features:

↳ Edge features are just like node features except it contains information about objects relationship/connection, not the object itself. E.g.



(1) Connection with bunch of attributes

(2) Connection with weights (E.g. how close 2 people are)