

Contents

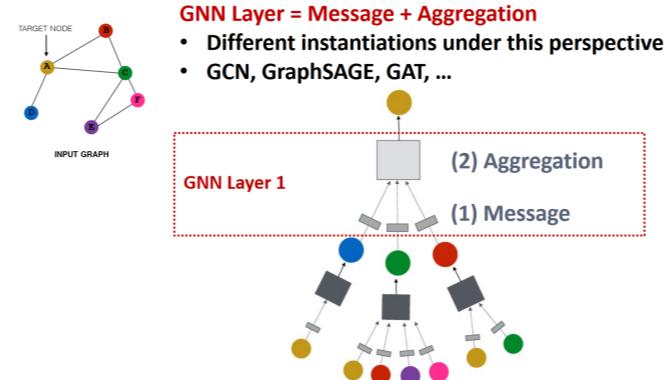
Overview.....	1
A general perspective on GNN.....	1
General GNN framework.....	1
1, 2: Message, aggregation.....	1
3: Layer connectivity	1
4: Graph augmentation.....	1
5: Learning objective	1
One GNN layer	1
Message computation.....	1
Issue & solutions – itself message lost.....	2
GCN baesd Summary	2
GCN as message transformation + aggregation.....	2
Details: normalise by node degree & sum	2
GraphSAGE.....	2
Different to GCN	2
Neighbour aggregation – mean, pool (MLP), LSTM	2
L2 Normalisation – layer norm 4 embedding	2
GAT	2
Not all node are important	2
Attention coefficient	3
Attention weighted sum aggregation; softmax [0,1].....	3
Attention Form.....	3
Multi-head attention.....	3
Benefit.....	3
In practice.....	3
Batch normalisation	3
Dropout	3
Activation (Non-linearity)	4
Stack GNN layers	4
Sequentially.....	4
The depth of GNN	4
Over-smoothing problem	4
Receptive field	4

Receptive field & over-smoothing.....	4
Solutions	4
get what's needed.....	4
Shallow GNN more expressive.....	4
1 - NN as aggregation	4
2 - Combine classic NN with GNN.....	4
Need many layers.....	4
1 - Add skip connections in GNNs	4
Idea of skip connection – mixture of models.....	5
GCN with skip connection.....	5
2- Directly skip to the last layer	5
Summary	5

Overview

A general perspective on GNN General GNN framework

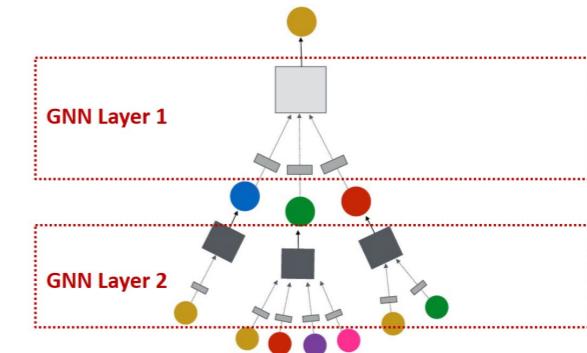
1, 2: Message, aggregation



3: Layer connectivity

Connect GNN layers into a GNN

- Stack layers sequentially
- Ways of adding skip connections



4: Graph augmentation

Idea: Raw input graph \neq computational graph

- Graph feature augmentation
- Graph structure augmentation

5: Learning objective

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

One GNN layer

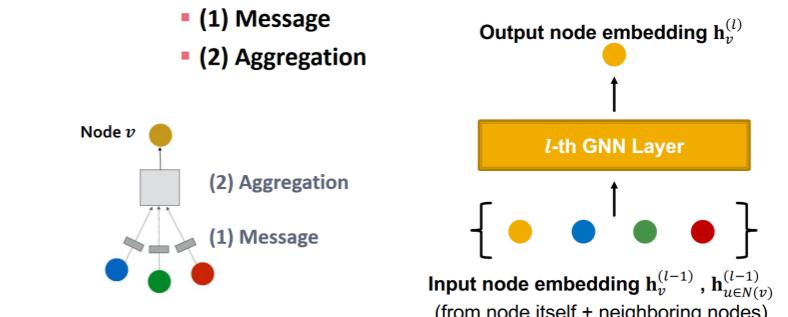
- The aggregation needs to be **order invariant**

Idea of a GNN Layer:

- Compress a set of vectors into a single vector

Two-step process:

- (1) Message
- (2) Aggregation



Message computation

- Message function:** $m_u^{(l)} = MSG^{(l)}(h_u^{(l-1)})$
 - Intuition:** Each node will create a message, which will be sent to other nodes later
 - Example:** A Linear layer $m_u^{(l)} = W^{(l)}h_u^{(l-1)}$
 - Multiply node features with weight matrix $W^{(l)}$

Issue & solutions – itself message lost

- Message.** Message from the node itself has different message computation compared to its neighbours. Denoted by B compared to W .
- Aggregations.** Aggregate itself embedding & concat with the neighbour aggregation
- Solution:** Include $\mathbf{h}_v^{(l-1)}$ when computing $\mathbf{h}_v^{(l)}$
 - (1) Message:** compute message from node v itself
 - Usually, a **different message computation** will be performed
 - (2) Aggregation:** After aggregating from neighbors, we can **aggregate the message from node v itself**
 - Via **concatenation or summation**

$$\mathbf{h}_v^{(l)} = \text{CONCAT} \left(\text{AGG} \left(\{\mathbf{m}_u^{(l)}, u \in N(v)\} \right), \mathbf{m}_v^{(l)} \right)$$

Then aggregate from node itself
First aggregate from neighbors

GCN based Summary

- Message is some transformation on the embedding
- Aggregate the message from the neighbour and itself
- Putting things together:**
 - (1) Message:** each node computes a message $\mathbf{m}_u^{(l)} = \text{MSG}^{(l)}(\mathbf{h}_u^{(l-1)}), u \in \{N(v) \cup v\}$
 - (2) Aggregation:** aggregate messages from neighbors $\mathbf{h}_v^{(l)} = \text{AGG}^{(l)}(\{\mathbf{m}_u^{(l)}, u \in N(v)\}, \mathbf{m}_v^{(l)})$
 - Nonlinearity (activation):** Adds expressiveness
 - Often written as $\sigma(\cdot)$: ReLU(\cdot), Sigmoid(\cdot), ...
 - Can be added to message or aggregation

GCN as message transformation + aggregation

(1) Graph Convolutional Networks (GCN)

$$\mathbf{h}_v^{(l)} = \sigma \left(\mathbf{W}^{(l)} \sum_{u \in N(v)} \frac{\mathbf{h}_u^{(l-1)}}{|N(v)|} \right)$$

Message transformation + aggregation.

$$\mathbf{h}_v^{(l)} = \sigma \left(\sum_{u \in N(v)} \frac{\mathbf{W}^{(l)} \mathbf{h}_u^{(l-1)}}{|N(v)|} \right)$$

Message
Aggregation

Details: normalise by node degree & sum

- Message:**
 - Each Neighbor: $\mathbf{m}_u^{(l)} = \frac{1}{|N(v)|} \mathbf{W}^{(l)} \mathbf{h}_u^{(l-1)}$ Normalized by node degree
(In the GCN paper they use a slightly different normalization)
- Aggregation:**
 - Sum over messages from neighbors, then apply activation
 - $\mathbf{h}_v^{(l)} = \sigma \left(\text{Sum} \left(\{\mathbf{m}_u^{(l)}, u \in N(v)\} \right) \right)$

In GCN graph is assumed to have self-edges that are included in the summation.

GraphSAGE

Different to GCN

- The aggregation has many choices
- Aggregation is based on the embedding over the computed message (some linear transformation of the embedding)

(2) GraphSAGE

$$\mathbf{h}_v^{(l)} = \sigma \left(\mathbf{W}^{(l)} \cdot \text{CONCAT} \left(\mathbf{h}_v^{(l-1)}, \text{AGG} \left(\{\mathbf{h}_u^{(l-1)}, u \in N(v)\} \right) \right) \right)$$

How to write this as Message + Aggregation?

- Message** is computed within the **AGG(\cdot)**
- Two-stage aggregation**
 - Stage 1:** Aggregate from node neighbors $\mathbf{h}_{N(v)}^{(l)} \leftarrow \text{AGG} \left(\{\mathbf{h}_u^{(l-1)}, u \in N(v)\} \right)$
 - Stage 2:** Further aggregate over the node itself $\mathbf{h}_v^{(l)} \leftarrow \sigma \left(\mathbf{W}^{(l)} \cdot \text{CONCAT}(\mathbf{h}_v^{(l-1)}, \mathbf{h}_{N(v)}^{(l)}) \right)$

Neighbour aggregation – mean, pool (MLP), LSTM

- LSTM is not **ordering invariance**
 - Need to permute the order
- Mean:** Take a weighted average of neighbors

$$\text{AGG} = \underbrace{\sum_{u \in N(v)} \frac{\mathbf{h}_u^{(l-1)}}{|N(v)|}}_{\text{Aggregation}} \quad \text{Message computation}$$
- Pool:** Transform neighbor vectors and apply symmetric vector function $\text{Mean}(\cdot)$ or $\text{Max}(\cdot)$

$$\text{AGG} = \underbrace{\text{Mean}([\text{MLP}(\mathbf{h}_u^{(l-1)}), \forall u \in N(v)])}_{\text{Aggregation}} \quad \text{Message computation}$$
- LSTM:** Apply LSTM to reshuffled of neighbors

$$\text{AGG} = \underbrace{\text{LSTM}([\mathbf{h}_u^{(l-1)}, \forall u \in \pi(N(v))])}_{\text{Aggregation}}$$

L2 Normalisation – layer norm 4 embedding

- Without [0, 1] normalisation, the embedding has different scale
- Optional:** Apply ℓ_2 normalization to $\mathbf{h}_v^{(l)}$ at every layer

$$\mathbf{h}_v^{(l)} \leftarrow \frac{\mathbf{h}_v^{(l)}}{\|\mathbf{h}_v^{(l)}\|_2} \quad \forall v \in V \text{ where } \|\mathbf{u}\|_2 = \sqrt{\sum_i u_i^2} \quad (\ell_2\text{-norm})$$
- Without ℓ_2 normalization, the embedding vectors have different scales (ℓ_2 -norm) for vectors
- In some cases (not always), normalization of embedding results in performance improvement
- After ℓ_2 normalization, all vectors will have the same ℓ_2 -norm

GAT

- Attention weight α_{vu} – the importance of a message coming from the node u (neighbour) for node v who is aggregating the message
- For GCN / GraphSAGE – neighbours of same importance to node v
- Intuition – allow the model to learn which part is important

(3) Graph Attention Networks

$$\mathbf{h}_v^{(l)} = \sigma \left(\sum_{u \in N(v)} \alpha_{vu} \mathbf{W}^{(l)} \mathbf{h}_u^{(l-1)} \right)$$

Attention weights

In GCN / GraphSAGE

- $\alpha_{vu} = \frac{1}{|N(v)|}$ is the **weighting factor (importance)** of node u 's message to node v
- $\Rightarrow \alpha_{vu}$ is defined **explicitly** based on the **structural properties of the graph** (node degree)
- \Rightarrow All neighbors $u \in N(v)$ are **equally important to node v**

Not all node are important

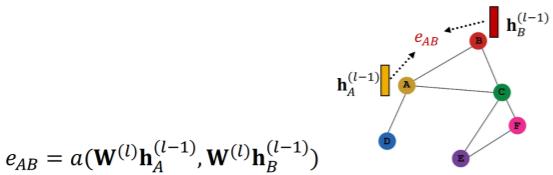
- Attention is inspired by cognitive attention.
- The **attention α_{vu}** focuses on the important parts of the input data and fades out the rest.
- Idea:** the NN should devote more computing power on that small but important part of the data.
- Which part of the data is more important depends on the context and is learned through training.

- Goal:** Specify **arbitrary importance** to different neighbors of each node in the graph

- Idea:** Compute embedding $\mathbf{h}_v^{(l)}$ of each node in the graph following an **attention strategy**:
 - Nodes attend over their neighborhoods' message
 - Implicitly specifying different weights to different nodes in a neighborhood

Attention coefficient

- Some function based on the previous layer embedding of the related two nodes
 - Let α_{vu} be computed as a byproduct of an **attention mechanism a** :
 - (1) Let a compute **attention coefficients e_{vu}** across pairs of nodes u, v based on their messages:
- $e_{vu} = a(\mathbf{W}^{(l)} \mathbf{h}_u^{(l-1)}, \mathbf{W}^{(l)} \mathbf{h}_v^{(l-1)})$
- e_{vu} indicates the importance of u 's message to node v



Attention weighted sum aggregation; softmax [0,1]

- Normalize e_{vu} into the final attention weight α_{vu}
 - Use the softmax function, so that $\sum_{u \in N(v)} \alpha_{vu} = 1$:
- $$\alpha_{vu} = \frac{\exp(e_{vu})}{\sum_{k \in N(v)} \exp(e_{vk})}$$

- Weighted sum based on the final attention weight α_{vu}

$$\mathbf{h}_v^{(l)} = \sigma(\sum_{u \in N(v)} \alpha_{vu} \mathbf{W}^{(l)} \mathbf{h}_u^{(l-1)})$$

Weighted sum using $\alpha_{AB}, \alpha_{AC}, \alpha_{AD}$:

$$\mathbf{h}_A^{(l)} = \sigma(\alpha_{AB} \mathbf{W}^{(l)} \mathbf{h}_B^{(l-1)} + \alpha_{AC} \mathbf{W}^{(l)} \mathbf{h}_C^{(l-1)} + \alpha_{AD} \mathbf{W}^{(l)} \mathbf{h}_D^{(l-1)})$$

Attention Form

- Learn jointly – attention weight & transform weight matrix

What is the form of attention mechanism a ?

- The approach is agnostic to the choice of a
 - E.g., use a simple single-layer neural network
 - a have trainable parameters (weights in the Linear layer)

$$\begin{array}{ccc} \text{Concatenate} & \xrightarrow{\text{Linear}} & e_{AB} \\ \mathbf{h}_A^{(l-1)} \quad \mathbf{h}_B^{(l-1)} & \xrightarrow{\text{Linear}} & \mathbf{e}_{AB} = a(\mathbf{W}^{(l)} \mathbf{h}_A^{(l-1)}, \mathbf{W}^{(l)} \mathbf{h}_B^{(l-1)}) \\ & & = \text{Linear}(\text{Concat}(\mathbf{W}^{(l)} \mathbf{h}_A^{(l-1)}, \mathbf{W}^{(l)} \mathbf{h}_B^{(l-1)})) \end{array}$$

- Parameters of a are trained jointly:

- Learn the parameters together with weight matrices (i.e., other parameter of the neural net $\mathbf{W}^{(l)}$) in an end-to-end fashion

Multi-head attention

- Stabilizes the learning process of attention mechanism
- Only 1 attention score may stuck at local minimum
 - Multi-head attention then aggregate them
 - The learning more stable
- When use, use each separately to calculate
- Weight initialisation should be different

- Create multiple attention scores (each replica with a different set of parameters):

$$\mathbf{h}_v^{(l)}[1] = \sigma(\sum_{u \in N(v)} \alpha_{vu}^1 \mathbf{W}^{(l)} \mathbf{h}_u^{(l-1)})$$

$$\mathbf{h}_v^{(l)}[2] = \sigma(\sum_{u \in N(v)} \alpha_{vu}^2 \mathbf{W}^{(l)} \mathbf{h}_u^{(l-1)})$$

$$\mathbf{h}_v^{(l)}[3] = \sigma(\sum_{u \in N(v)} \alpha_{vu}^3 \mathbf{W}^{(l)} \mathbf{h}_u^{(l-1)})$$
- Outputs are aggregated:
 - By concatenation or summation
 - $\mathbf{h}_v^{(l)} = \text{AGG}(\mathbf{h}_v^{(l)}[1], \mathbf{h}_v^{(l)}[2], \mathbf{h}_v^{(l)}[3])$

Benefit

- Key benefit:** Allows for (implicitly) specifying **different importance values (α_{vu}) to different neighbors**
- Computationally efficient:**
 - Computation of attentional coefficients can be parallelized across all edges of the graph
 - Aggregation may be parallelized across all nodes
- Storage efficient:**
 - Sparse matrix operations do not require more than $O(V + E)$ entries to be stored
 - Fixed number of parameters, irrespective of graph size
- Localized:**
 - Only attends over local network neighborhoods
- Inductive capability:**
 - It is a shared edge-wise mechanism
 - It does not depend on the global graph structure

In practice

Batch Normalization:

Stabilize neural network training

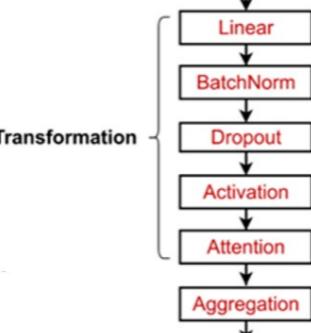
Dropout:

Prevent overfitting

Attention/Gating:

Control the importance of a message

A suggested GNN Layer



Batch normalisation

- Goal:** Stabilize neural networks training
- Idea:** Given a batch of inputs (node embeddings)
 - Re-center the node embeddings into zero mean
 - (for each feature in the col)
 - Re-scale the variance into unit variance

Input: $\mathbf{X} \in \mathbb{R}^{N \times D}$
N node embeddings

Trainable Parameters:
 $\gamma, \beta \in \mathbb{R}^D$

Output: $\mathbf{Y} \in \mathbb{R}^{N \times D}$
Normalized node embeddings

$$\mu_j = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_{i,j}$$

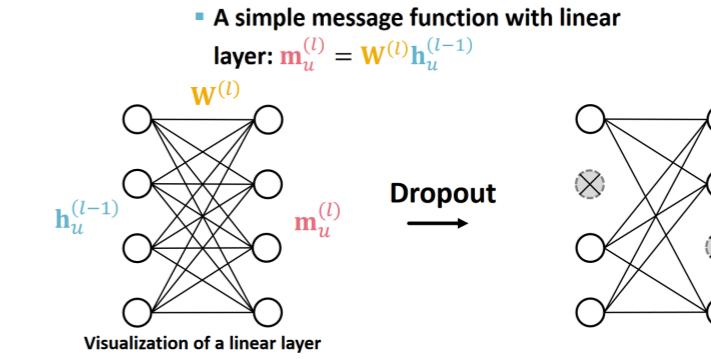
$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_{i,j} - \mu_j)^2$$

$$\hat{\mathbf{x}}_{i,j} = \frac{\mathbf{x}_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

$$\mathbf{y}_{i,j} = \gamma_j \hat{\mathbf{x}}_{i,j} + \beta_j$$

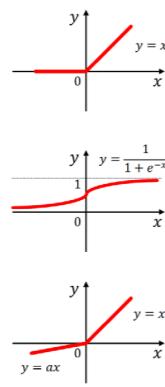
Dropout

- Goal:** Regularize a neural net to prevent overfitting
- During **training**: with some probability p , randomly set neurons to zero (turn off)
- During **testing**: Use all the neurons for computation
- In GNN, Dropout is applied to the linear layer in the message function



Activation (Non-linearity)

- Apply activation to i -th dimension of embedding x
- Rectified linear unit (ReLU)**
 $\text{ReLU}(x_i) = \max(x_i, 0)$
 Most commonly used
- Sigmoid**
 $\sigma(x_i) = \frac{1}{1 + e^{-x_i}}$
 Used only when you want to restrict the range of your embeddings
- Parametric ReLU**
 $\text{PReLU}(x_i) = \max(x_i, 0) + a_i \min(x_i, 0)$
 a_i is a trainable parameter
 Empirically performs better than ReLU

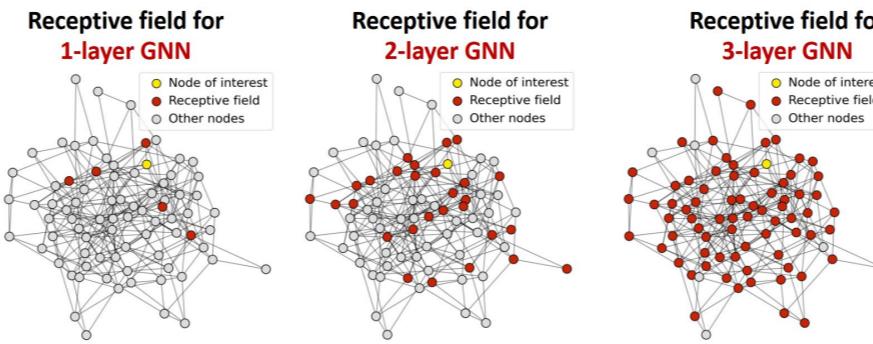


each node are similar

Receptive field

Receptive field:

- the set of nodes that determine the embedding of a node of interest
 - In a K -layer GNN, each node has a receptive field of K -hop neighborhood
- 2 layers
- Neighbour of neighbours
 - 1 hop or 2 hops neighbourhood
 - The shared neighbours (in red) between the two yellow nodes quickly grows when we increase the number of hops (num of GNN layers)



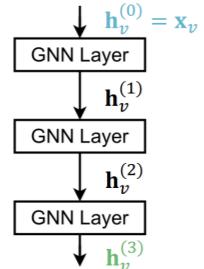
Stack GNN layers

How to connect GNN layers into a GNN?

- Stack layers **sequentially**
- Ways of adding **skip connections**

Sequentially

- The standard way:** Stack GNN layers sequentially
- Input:** Initial raw node feature x_v
- Output:** Node embeddings $h_v^{(L)}$ after L GNN layers



The depth of GNN

- Tells how many hops to aggregate
- Different to CNN
 - Not about the expressiveness of the model

Over-smoothing problem

- Issue of **stacking many** GNN layers
- all the node embeddings converge to the **same value**
 - all the node embeddings converge to the same value
 - bad because we want to use node embeddings to differentiate nodes
- reason.**
 - The receptive field is too large and so the info collected for

each node are similar

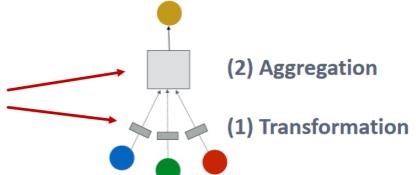
Shallow GNN more expressive

- Expressive power for one layer

1 - NN as aggregation

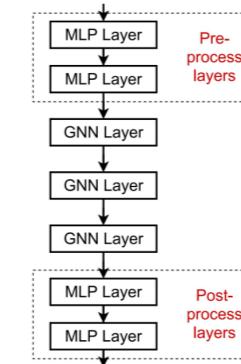
- We can make **aggregation / transformation** become a **deep neural network!**

If needed, each box could include a **3-layer MLP**



2 - Combine classic NN with GNN

- E.g., we can add **MLP layers** (applied to each node) before and after GNN layers, as **pre-process layers** and **post-process layers**



Pre-processing layers: Important when encoding node features is necessary.
E.g., when nodes represent images/text

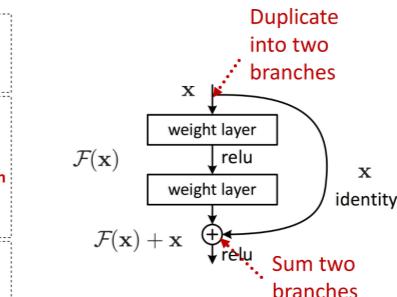
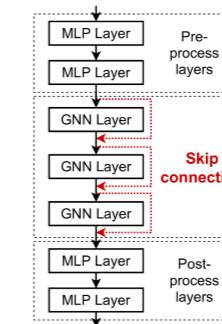
Post-processing layers: Important when reasoning / transformation over node embeddings are needed
E.g., graph classification, knowledge graphs

In practice, adding these layers works great!

Need many layers

1 - Add skip connections in GNNs

- Observation from over-smoothing: Node embeddings in **earlier GNN layers** can sometimes better differentiate node
- Solution:** We can **increase the impact** of earlier layers on the final node embeddings, by **adding shortcuts** in GNN

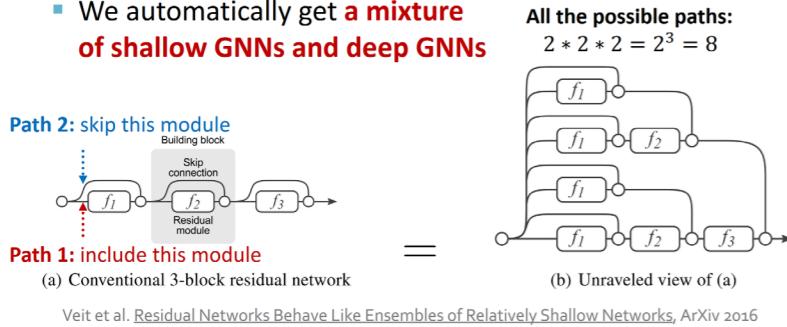


Idea of skip connections:
Before adding shortcuts:
 $F(x)$
After adding shortcuts:
 $F(x) + x$

Jure Leskovec, Stanford CS224W: Machine Learning with Graphs, <http://cs224w.stanford.edu>

Idea of skip connection – mixture of models

- **Why do skip connections work?**
 - **Intuition:** Skip connections create a **mixture of models**
 - N skip connections $\rightarrow 2^N$ possible paths
 - Each path could have up to N modules
 - We automatically get a **mixture of shallow GNNs and deep GNNs**

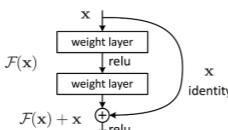


GCN with skip connection

A standard GCN layer

$$\mathbf{h}_v^{(l)} = \sigma \left(\sum_{u \in N(v)} \mathbf{W}^{(l)} \frac{\mathbf{h}_u^{(l-1)}}{|N(v)|} \right)$$

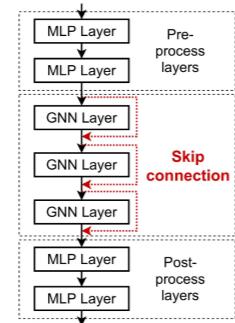
This is our $F(\mathbf{x})$



A GCN layer with skip connection

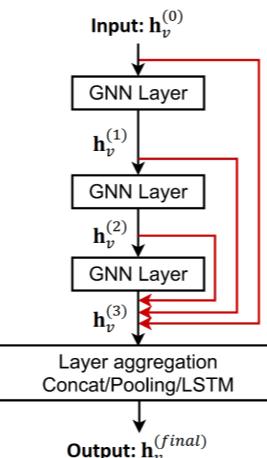
$$\mathbf{h}_v^{(l)} = \sigma \left(\sum_{u \in N(v)} \mathbf{W}^{(l)} \frac{\mathbf{h}_u^{(l-1)}}{|N(v)|} + \mathbf{h}_v^{(l-1)} \right)$$

$\mathbf{F}(\mathbf{x})$ + \mathbf{x}



2- Directly skip to the last layer

- The **final layer** directly aggregates from all the node embeddings in the **previous layers**
 - E.g. add info from early layer e.g. 1-hop neighbour or deeper ones e.g. 4-hop neighbour
- Xu et al. Representation learning on graphs with jumping knowledge networks, ICML 2018



AGG

- **Mean:** Take a weighted average of neighbors

$$\text{AGG} = \sum_{u \in N(v)} \frac{\mathbf{h}_u^{(l-1)}}{|N(v)|}$$

Aggregation Message computation

- **Pool:** Transform neighbor vectors and apply symmetric vector function $\text{Mean}(\cdot)$ or $\text{Max}(\cdot)$

$$\text{AGG} = \text{Mean}(\{\text{MLP}(\mathbf{h}_u^{(l-1)}) \mid \forall u \in N(v)\})$$

Aggregation Message computation

- **LSTM:** Apply LSTM to reshuffled of neighbors

$$\text{AGG} = \text{LSTM}([\mathbf{h}_u^{(l-1)} \mid \forall u \in \pi(N(v))])$$

Aggregation

GAT

$$\mathbf{h}_v^{(l)} = \sigma \left(\sum_{u \in N(v)} \alpha_{vu} \mathbf{W}^{(l)} \mathbf{h}_u^{(l-1)} \right)$$

Attention weights

Attention coefficient

- **Multi-head**
 - Instead of one, use many independent heads to get different attention coefficient measuring the importances of neighbours to the target node to avoid local minimum

$$e_{vu} = a(\mathbf{W}^{(l)} \mathbf{h}_u^{(l-1)}, \mathbf{W}^{(l)} \mathbf{h}_v^{(l-1)})$$

Over-smoothing problems

- collect too many information from neighbours and if all nodes do so, then very similar
 - Too many layers
 - layer in GNN refers to receptive field i.e. neighbour of a node (info to be collected)
- E.g. if receptive field is large enough to cover the whole graph

Receptive field & over-smoothing

- Treat increase in #GNN as an increase in receptive field

Solution – increase expressive power of GNN

1. NN as aggregation
2. Skip connection
 - a. Skip to last layer

GraphSage

$$\mathbf{h}_v^{(l)} = \sigma \left(\mathbf{W}^{(l)} \cdot \text{CONCAT} \left(\mathbf{h}_v^{(l-1)}, \text{AGG} \left(\{\mathbf{h}_u^{(l-1)} \mid \forall u \in N(v)\} \right) \right) \right)$$

Stage 1: Aggregate from node neighbors

$$\mathbf{h}_{N(v)}^{(l)} \leftarrow \text{AGG} \left(\{\mathbf{h}_u^{(l-1)} \mid \forall u \in N(v)\} \right)$$

Stage 2: Further aggregate over the node itself

$$\mathbf{h}_v^{(l)} \leftarrow \sigma \left(\mathbf{W}^{(l)} \cdot \text{CONCAT} (\mathbf{h}_v^{(l-1)}, \mathbf{h}_{N(v)}^{(l)}) \right)$$