

MESSAGE-PASSING : Information is propagated between nodes of a graph along its edges.

A GNN layer can be thought of as a message-passing setup where nodes update their embeddings by aggregating messages from their direct neighbours.

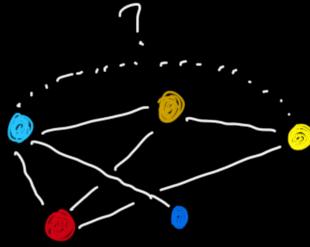
Various GNN models differ in how nodes aggregate the messages from their neighbours.

(I) GCN (Kipf and Welling) Element-wise mean of the embeddings of the neighbours.

(II) GraphSAGE (Hamilton et al) Element-wise Maximum of the embeddings of the nbrs.

(III) Graph Attention Network (Velickovic et al) (GAT)
weighted sum of the embeddings of the neighbours
where the weights for each neighbour are computed
via an attention b/w the current node and the
neighbour.

Often require to model the interaction b/w two non-connected nodes.



MP-GNNs achieve this by using multiple layers.



OVER-SQUASHING

Growing amount of information is squashed into fixed-sized vectors.

Message-passing fails to propagate information effectively over the graph.

⇒ MP-GNNs are not adept at capturing long-range dependencies.

FURTHERMORE

OVER-SMOOTHING Embeddings of the nodes become indistinguishable with increasing the depth of the network.



GNN Layers

Receptive field of a node grows exponentially with the radius of the neighborhood.

Traffic Analogy (Due to Dominique Beaini)

Over-Sampling \rightarrow Stuck in traffic

Over-Smoothing \rightarrow Lost in traffic

\Rightarrow Also relates to the expressivity of GNNs.
(in terms of WL-test)

GRAPH TRANSFORMERS

attempt to alleviate these problems of MP-GNNs.



Instead of passing messages b/w neighbours, the idea is to directly model the dependence b/w any two nodes in the graph.

As opposed to the local attention over neighbours in GAT, graph transformers seek to apply global attention over the whole graph.

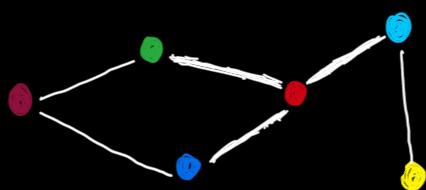
\Rightarrow Both connected and non-connected nodes can contribute to the node embedding.

Problem: Standard transformer with global attention is oblivious to the graph structure.

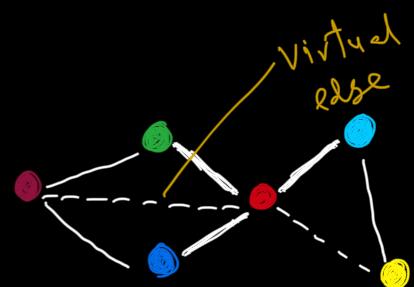
↓

Challenge is how to incorporate the structural information of the graph.

Different Graph transformer models differ in how they leverage graph structure along with the global attention.



MP-GNN



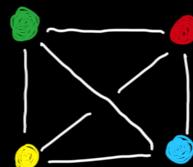
Graph Transformer

NEXT :

- Spectral Attn Network
- GraphiT
- Graphomer
- SignNet and BasisNet
- Learnable structural and positional Encoding.

TRANSFORMERS (NLP) treat a sentence as a fully-connected graph.

This is a sentence \leftrightarrow



Each word is connected to every other word.

Positional Encoding is chosen to account for the relative positions of the words within a sentence.

Transformer Input is Word Embedding \oplus Positional Encoding

How to extend the transformers to arbitrary graphs?

\downarrow
use Positional Encoding (PE) for nodes.

Challenging to design Positional Encoding for graphs since there is NO canonical ordering on graphs.

Indeed, GNNs are designed to learn node embeddings that are invariant to the node positions.
(Permutation Invariance)

DESIDERATA for PE :

(I) PE should be chosen such that ^{the} nearby nodes have similar PE_x and the nodes that are far apart have dissimilar PE_x .

(II) PE of a node should be unique.

↓
But graph structure imposes Symmetries.

E.g. if nodes i and j are structurally symmetric (interchangeable)
then their PE_x are also interchangeable.

$$i \leftrightarrow j \Rightarrow PE(i) \leftrightarrow PE(j)$$

⇒ PE is arbitrary upto symmetries of the graph.

(L) NEED to design architectures that can deal with these ambiguities.

Refined (II) The amount of ambiguities in PE should be as little as possible.

Simplest possible choice for PE is to assign an ordering to the nodes. $\forall v \in V, PE(v) = \text{Index}(v)$

However, there are $N!$ possible orderings. $N = |V|$



Number of Nodes

During training the orderings are uniformly sampled from $N!$ possible choices

such that the network is independent of the choice of the ordering.

(Murphy et al 2019)

Proposal : use graph Laplacian eigen vectors as PE.

(Dwivedi et al
2021)

- Provide an embedding of the nodes into the Euclidean space.
- A meaningful coordinate system.
- Cognizant of the global graph structure.
- Encode distance aware information
 - [Nearby nodes have similar PE,
Far-apart nodes have different PE]

⇒ Laplacian eigenvectors can be interpreted as a GENERALIZATION of the PE of the NLP transformer.

Sentence graph Laplacian = Discrete Line

\downarrow
eigenvectors
Sine and cosine

$$\forall v \in V, \quad PE(v) = (\underbrace{\lambda_1, \lambda_2, \dots, \lambda_k}_{k\text{-smallest non-trivial eigenvectors}}) \in \mathbb{R}^k$$

of the Laplacian.

$$\Delta = I - D^{-1/2} A D^{-1/2} \xrightarrow{\text{EIGEN DECOMPOSITION}} U^T \Lambda U$$

↗ ↘
 Laplacian Degree matrix Adjacency matrix eigenvec.
 ↘ ↗
 eigenvalues

Q: Are these Positional Encodings Unique?

NO

Laplacian eigenvectors are defined upto sign.

$\pm 1^y$ symmetry.

↓
Network needs to deal with this ambiguity.

During training the sign of the eigenvectors are randomly flipped so as to have the model invariant

[Dwivedi et al 2021]

+ sign.



Possible number of flips 2^k .

■ Laplacian PE reduces the sample space of the
~~~~~ node-ordering PE.



Invariant to node-reparametrization.

But defined upto the sign.

Choose  $k \ll N \implies 2^k \ll N!$

Laplacian PE has been used in Various Graph Transformer architectures.

An added benefit is that Laplacian eigenvectors can be pre-computed.

# Laplacian Positional Encoding

PE in the NLP transformer  $\rightarrow$  Sine and Cosine functions  
(Vaswani et al)

$\downarrow$   
CAN'T be defined for arbitrary graphs  
Since there is no notion  
of position along an axis.  
(time-steps for words in  
a sentence)

Extension of Sinusoidal fns to arbitrary graphs



Eigenvectors of the graph Laplacian.

- In Euclidean space, Laplacian eigenfunctions are Sine/Cosine with squared frequencies as r-values.

Dwivedi et al proposed using  $k$  Laplacian eigenvectors as PE.  
(fixed #)

(A graph  $G$  with  $N$  nodes can have at most  $N$  linearly independent eigenvectors)

Now the graphs in a dataset can have varying numbers of nodes.



Bottleneck when there is a large disparity in the number of nodes across graphs.  
i.e. the smallest graphs have significantly fewer nodes than the largest graphs.



A small fraction of the eigenvectors will be used for large graphs.



Loss of information!

- This motivates the need for PE of dimension  $k$  which DOES NOT depend on the number of eigenvectors in the graph



Spectral Attention Network

(Krauzer et al NeurIPS 2021)

↓  
utilize full spectrum of the Laplacian spectrum to learn the position of each node in the graph.

[LEARNABLE POSITIONAL ENCODING]  
(LPE)

Idea:

Can use the Laplacian eigenvectors as initialization of the trainable PE.

(I) Calculate Laplacian eigenvectors of the graph.  
initial static PE

(II) use a suitable architecture to process the LapPE  
to "learn" PE for the nodes.

→ Can be a Neural Network

→ Can be learnt End-to-End with the Graph transformer

OR Can be an independent unsupervised or self-supervised component.

⇒ Krueger et al proposed a transformer type architecture to process initial Laplacian PE into learnable PE which are independent of the number of Laplacian eigenvectors.

Transformers are most effective on sequential data  
(with PE)

Meta Question: How to think of Laplacian eigenvectors as sequences?

i.e. What is the appropriate PE for Laplacian Eigenvectors.

Fourier transform maps a function from space (or time) domain to spatial (or temporal) frequency domain.

In Euclidean space, for a fn

$$\hat{f}(t) = \int_{\mathbb{R}^n} f(x) e^{-2\pi i t \cdot x} dx$$

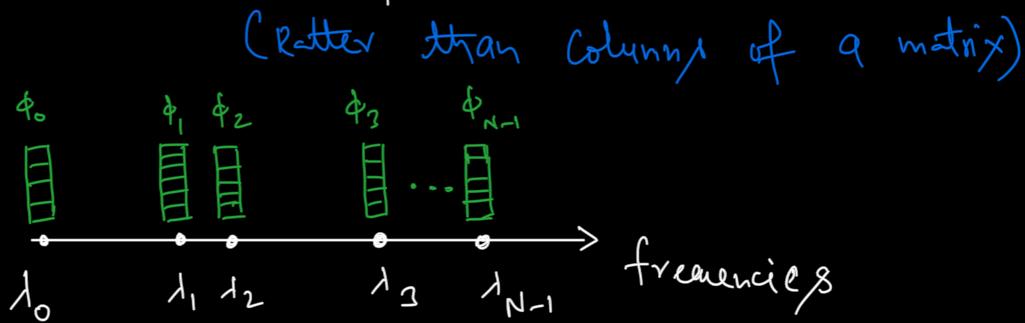
$\downarrow$  Euler's formula  
 $\cos(2\pi t \cdot x) - i \sin(2\pi t \cdot x)$

For graphs, the Fourier transform of a fn is defined as

$$F(f)(\lambda_i) = \langle f, \phi_i \rangle$$

$\downarrow$  eigenvalue of the Laplacian       $\downarrow$  eigenvector of the Laplacian  
 the Laplacian corresponding to  $\lambda_i$ .

$\Rightarrow$  Laplacian eigenvectors are best viewed as vectors positioned on the axis of eigenvalues.



$\Rightarrow$  We can think of Laplacian eigenvectors as sequences with the corresponding eigenvalues as fine-steps (Non-uniform)

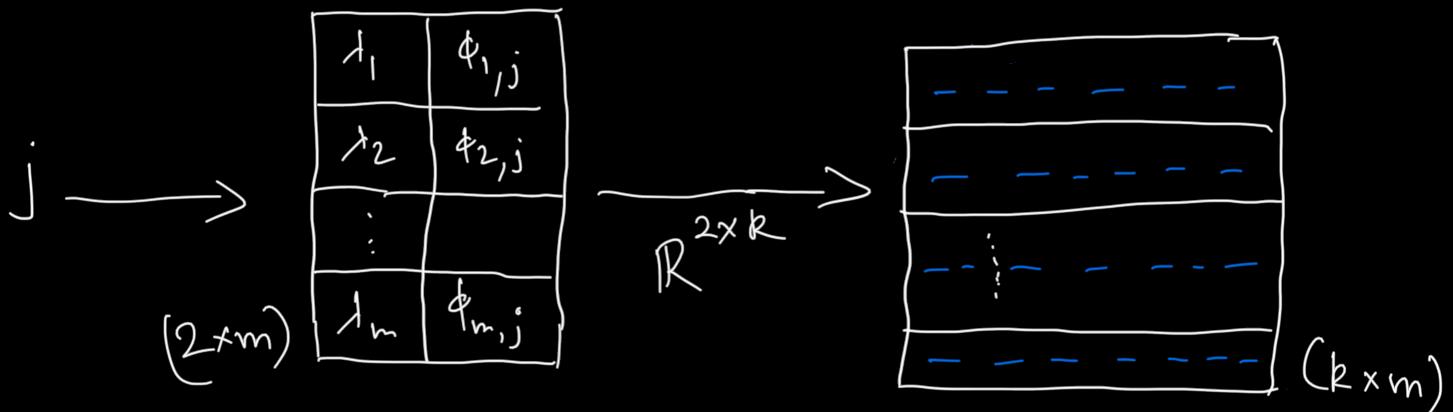
## LPE transformer à la Kreuzer et al.

- (I) Create Embedding matrix for each node  $j$  by  
 $(2 \times m)$   
 Concatinating the  $m$ -smallest eigenvalues  
 with their corresponding eigenvectors.

$$j \rightarrow \begin{array}{|c|c|} \hline \lambda_1 & \phi_{1,j} \\ \hline \lambda_2 & \phi_{2,j} \\ \hline \vdots & \vdots \\ \hline \lambda_m & \phi_{m,j} \\ \hline \end{array}$$

- =  $m$  is a hyperparameter similar to length of sentence.
- = For graphs where  $m > \# \text{nodes}$ , a padding is added.

- (II) A Linear projection layer is applied to map the initial embeddings to  $k$ -dimensional Embeddings. ( $k$  is a hyperparam)



- (III) A transformer encoder then computes self attention over the sequence of length  $m$  and hidden dimension  $k$ .

- (IV) Sum-pooling along the  $m$ -axis produces an embedding of fixed  $k$ -dimension.  
 $\Downarrow$   
 Learned Positional Embedding of dimension  $k$ .

Inspired by Physics

(Pseudo) Inverse of the Laplacian represents electrostatic potential of a charge.  
Green's function

Let  $\Delta = D^{-1/2} L D^{-1/2}$  be the normalized, symmetric Laplacian.

$\lambda_i^{\hat{y}}$  and  $\phi_i^{\hat{y}}$

↓  
evalues and ↓ eigenvectors of  $\Delta$ .

$$\Delta^{-1} := G(j_1, j_2) = \underbrace{d_{j_1}^{1/2} d_{j_2}^{-1/2}}_{\downarrow \text{electric potential}} \underset{i > 0}{\lesssim} \left( \frac{\phi_{i,j_1} \phi_{i,j_2}}{\lambda_i} \right)^2$$

between nodes  $j_1$  &  $j_2$ .

$\Rightarrow$  product of the Laplacian eigenvectors is proportional to the electrostatic interaction between 2 nodes.

RECALL: we can build node positional encodings (PE) using the eigenvectors of the graph Laplacian.

## Let's build Graph Transformers

Let  $G = (V, E)$  be a graph.

Nodes  $\hookrightarrow$  Edges

if node  $i \in V$ ,  $\alpha_i \in \mathbb{R}^{d_n}$  Node features

if edge  $(i, j) \in E$ ,  $\beta_{ij} \in \mathbb{R}^{d_e}$  Edge features

if node  $i \in V$ ,  $d_i \in \mathbb{R}^k$  Node positional Encoding.

The first step in the GT architecture is the linear projections of the node feats, edge feats, and the node PE.

$$\hat{h}_i^\circ = A^\circ \alpha_i + a^\circ, \quad A^\circ \in \mathbb{R}^{d \times d_n}, \quad a^\circ \in \mathbb{R}^d$$

$$e_{ij}^\circ = B^\circ \beta_{ij} + b^\circ, \quad B^\circ \in \mathbb{R}^{d \times d_e}, \quad b^\circ \in \mathbb{R}^d$$

$$d_i^\circ = C^\circ d_i + c^\circ, \quad C^\circ \in \mathbb{R}^{d \times k}, \quad c^\circ \in \mathbb{R}^d$$

Input Node Embedding  $h_i^\circ = \hat{h}_i^\circ + d_i^\circ$   
 (Same as NLP transformers)

## $(l+1)$ -th Layer of the GT (Dwivedi et al)

H node  $i$

$$\hat{h}_i^{l+1} = \odot_h^l \underset{k=1}{\underset{\text{Concat}}{\parallel}} \left( \sum_{j \in \mathcal{N}_i} w_{ij}^{k,l} v^{k,l} h_j^l \right)^{\#(\text{heads})}$$

$\downarrow$   
Sum over Neighbours

$$\text{Softmax}_j \left( \frac{Q^{k,l} h_i^l \cdot K^{k,l} h_j^l}{\sqrt{d_k}} \right)$$

$$h_i^l \rightarrow Q^{k,l} h_i^l \quad \text{QUERY}$$

$$h_j^l \rightarrow K^{k,l} h_j^l \quad \text{KEY}$$

$$h_j^l \rightarrow V^{k,l} h_j^l \quad \text{VALUE}$$

$$\begin{matrix} Q^{k,l} \\ K^{k,l} \\ V^{k,l} \end{matrix} \} \in \mathbb{R}^{d_k \times d_l}$$

$\Rightarrow$  Query attends to the values based on its correlations (Neighbours) with the corresponding key.  
 (Current Node)

This is direct application of the Scaled-dot Product Attn of the transformer to graphs.  $\text{Softmax} \left( \frac{K Q^T}{\sqrt{d_k}} \right) V$

Graph Attention Network (GAT) used additive attention

(Velickovic et al)

$$w_{ij} = \text{softmax}_j \left( \text{MLP}(W h_i \oplus W h_j) \right)$$

⇒ The attention outputs across heads are concatenated and then linearly projected to obtain the  $(l+1)$ -th layer node embedding.

The node embeddings pass through Norm and feed-fwd layers with residual connections before fed to the task-specific network.

$$\hat{h}_i^{l+1} = \text{Norm} \left( h_i^l + \hat{h}_i^{l+1} \right)$$

$$\hat{h}_i^{l+1} = w_2^l \text{ReLU} \left( w_1^l \hat{h}_i^{l+1} \right)$$

$$h_i^{l+1} = \text{Norm} \left( \hat{h}_i^{l+1} + \hat{h}_i^{l+1} \right)$$

$$\text{output} = \text{MLP}(h_i^L) \rightarrow \text{embedding after } L\text{-layers.}$$

GAT Layer with Edge features

Motivated by the better utilization of rich information present in several datasets in the edge features,

Dwivedi et al proposed a modification of the GAT layer in the presence of edge features.

$$\hat{\omega}_{ij}^{k,l} = \left( \frac{Q^{k,l} h_i^l \cdot K^{k,l} h_j^l}{\sqrt{d_k}} \right) \cdot E^{k,l} e_{ij}^l$$

edge features  
 $\rightarrow \mathbb{R}^{d_k \times d_l}$

$$\omega_{ij}^{k,l} = \text{softmax}_j(\hat{\omega}_{ij}^{k,l})$$

$$\hat{h}_i^{l+1} = \odot_h^l \left| \begin{array}{l} \parallel \\ \leq \sum_{j \in N_i} \omega_{ij}^{k,l} V^{k,l} h_j^l \end{array} \right|_{k=1}^K$$

Edge features can be updated as

$$\hat{e}_{ij}^{l+1} = \odot_e^l \left| \begin{array}{l} \parallel \\ \left( \hat{\omega}_{ij}^{k,l} \right) \end{array} \right|_{k=1}^K$$

Think of  $\hat{\omega}_{ij}^{k,l} = Q^{k,l} h_i^l \cdot K^{k,l} h_j^l$  as the intermediate attention score quantifying the information about the edge  $(i,j)$ .

When additional edge info in the form of edge features is available,  $\hat{\omega}_{ij}^{k,l}$  is modified as —

$$\hat{\omega}_{ij}^{k,l} = \hat{\omega}_{ij}^{k,l} * E^{k,l} e_{ij}^l$$



## Graph Transformer a la Srivastava et al

- Local attention over the neighbourhood of each node i.e. the attn depends on the neighbourhood connectivity
- The node PE is given by the eigenvectors of the graph Laplacian.
- Architecture can be extended to include edge features.



- Still local attention
- LapPE can result in bottlenecks.



## Spectral Attention Network :

- Global full-graph Attention
- Learnable Positional Encodings (LPE)

Full-graph Attn while preserving the local connectivity of the nodes in the graph.

use two sets of attention mechanisms

- one for nodes connected by real edges
- one for nodes connected by virtual edges  
(Any pair of nodes)

Similar to the GT architecture, the  $(l+1)$ -th layer of SAN is defined as

$$\hat{h}_i^{l+1} = \odot_h^l \left( \sum_{j \in V} w_{ij}^{k,l} v^{k,l} h_j^l \right)$$

$\hookrightarrow$  over all Nodes

where  $w_{ij}^{k,l} = \begin{cases} \frac{1}{1+\gamma} \text{softmax}(\hat{w}_{ij}^{k,l}) & \text{if } (i,j) \in \mathcal{E} \\ \frac{\gamma}{1+\gamma} \text{softmax}(\hat{w}_{ij}^{k,l}) & \text{if } (i,j) \notin \mathcal{E} \end{cases}$

$\gamma \in \mathbb{R}^+$  is a hyperparameter that controls the amount of bias towards the full-graph attention.

SAN uses the edge feature enhanced version of the GT attention as follows -

$$\hat{w}_{ij}^{k,l} = \begin{cases} \frac{Q^{1,k,l} h_i^l \cdot K^{1,k,l} h_j^l \cdot E^{1,k,l} e_{ij}}{\sqrt{d_k}} & \text{if } (i,j) \in \mathcal{E} \\ \frac{Q^{2,k,l} h_i^l \cdot K^{2,k,l} h_j^l \cdot E^{2,k,l} e_{ij}}{\sqrt{d_k}} & \text{if } (i,j) \notin \mathcal{E} \end{cases}$$

Similar to GT, we have Norm and Residual Connection.

$$\hat{h}_i^{l+1} = \text{Norm}(\hat{h}_i^l + h_i^l)$$

$$\hat{h}_i^{l+1} = w_2^l \text{ReLU}(w_1^l \hat{h}_i^l)$$

$$h_i^{l+1} = \text{Norm}(\hat{h}_i^{l+1} + \hat{h}_i^{l+1})$$

Leveraging eigenvectors of the graph Laplacian for building node positional encodings (PE) has shown promise to craft expressive graph Transformer architectures.

↓  
(cf Dwivedi et al Graph Transformer  
Kremer et al Spectral Att<sub>n</sub> Network)

Limitation is that the Laplacian eigenvectors are defined upto the sign.

i.e. if  $v$  is an eigenvector, then so is  $-v$

↓

Need to have Sign-invariant PE.

Prior works dealt with the sign ambiguity by randomly flipping the sign of the eigenvectors during training.

↓

SignNet a new neural architecture that is invariant to sign flips.  
(Lim et al)

More generally, for any basis Symmetry i.e. infinitely many  
↓ choices for basis.

BasisNet - Invariant to basis Symmetry.

Idea : Process each of the eigenspaces individually, and then aggregate. → Spanned by the diff. choices of eigenvectors.

Let  $\{v_1, v_2, \dots, v_R\}$  be the Laplacian eigenvectors.

$\{\downarrow v_1, -v_2, \dots, -v_k\}$  are also the Laplacian eigenvectors.

Justification for this design (in 1d)

$\Leftrightarrow$  A continuous fn  $h: \mathbb{R}^n \rightarrow \mathbb{R}^S$  is sign-invariant if and only if  $h(x) = \phi(x) + \phi(-x)$  for some continuous fn  $\phi: \mathbb{R}^n \rightarrow \mathbb{R}^S$ .

Further,  $h$  is permutation-equivariant if & only if  $\exists \phi$  as above and  $\phi$  is permutation equivariant.  
 $Pf(x) = f(Px)$   
 $P \in \mathbb{R}^{n \times n}$  permutation matrix

$\Leftrightarrow$  without details, BasNet can be defined as -

$$f(v_1, v_2, \dots, v_k) = \phi \left( \left[ \phi_{d_i}(v_i v_i^T) \right]_{i=1, \dots, k} \right)$$

Order 2 invariant GNN  
MLP, DeepSets, Transformer.



## Graphomer

Ying et al (NeurIPS 2021) address the question of efficacy of the Transformer architecture in graph representation learning.

Recall, the key to adaptability of the transformer architecture to graphs is proper incorporation of structural information of graphs.

Self-attention operates over a pair of entities (nodes) and encodes the semantic similarity b/w them.



Agnostic to the relational information.

## Two Key Ideas

(1) Different nodes may have different importance e.g. celebrities tend to have more influence.

↓  
Encode node centrality explicitly into the model

↓  
Add learnable vector encoding to each node which is dependent on the degree of the node.

for each node  $v_i \in V$ ,

$$h_i^{(0)} = x_i + z_{\deg^-(v_i)}^- + z_{\deg^+(v_i)}^+$$

where  $z^-$ ,  $z^+ \in \mathbb{R}^d$  are learnable embedding vectors specified by the indegree ( $\deg^-$ ) and outdegree ( $\deg^+$ ) of the node.

(II) Different node pairs have different structural properties owing to the relational info.



Encode relational properties of node-pairs into the model.



Add learnable encoding for each node-pair which is dependent on the relational structure

Spatial Encoding  $\phi: V \times V \rightarrow \mathbb{R}$

$\phi(v_i, v_j)$  defines the connectivity b/w the nodes i and j.

$$\phi(v_i, v_j) := \begin{cases} \text{Shortest Path Distance b/w } v_i \text{ and } v_j \text{ if connected} \\ -1, \text{ otherwise.} \end{cases}$$

(II) Various graphs come with edge features which carry important structural information (e.g. in molecular graphs, atom pairs may have features describing properties of the chemical bonds between them).



Encode edge features into the model.



Add learnable edge embeddings which depend on the edge features as well as the graph structure.

A pair of nodes  $(v_i, v_j)$

Let  $SP_{ij} = (e_1, e_2, \dots, e_N)$  be the shortest path b/w them.

and  $w_n^E$  be the learnable vector encoding of the edge  $e_n$ .

then  $c_{ij} = \underbrace{\frac{1}{N} \sum_{n=1}^N}_{\text{Average along the shortest path.}} \underbrace{x_{e_n} (\underbrace{w_n^E}_{\text{original edge feature}})^T}_{\text{Dot-product of the given and learnable encoding.}}$

Finally, the Self-attention of the Graphmer is defined as -

$$A_{ij} = \underbrace{\frac{(h_i w_Q) (h_j w_K)^T}{\sqrt{d}}}_{\text{Self-attn of the standard transform}} + \underbrace{b(\phi(v_i, v_j))}_{\text{bias-term (fn of } \phi(v_i, v_j)\text{)}} + \underbrace{c_{ij}}_{\text{edge facts bias-term}}$$