

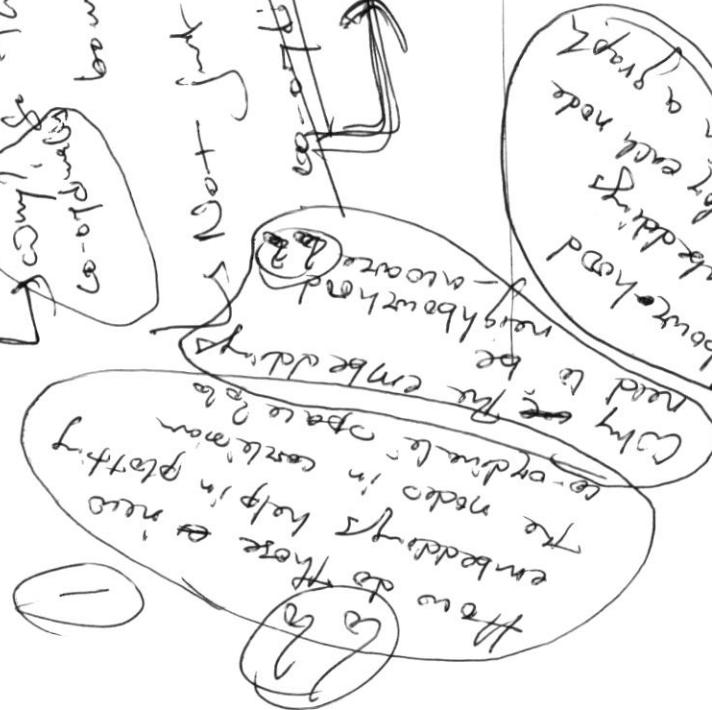
① 1.1

Vid 0 1: → Neural Message Passing

learn meaningful co-ordinates  
for each node in the graph  
( $\rightarrow$  k-dimensional embeddings) for each node to encode a node's features and  
not the ~~whole~~ division

graph  
containing nodes for classification  
co-ordinates for node neighbour  
embeddings  
Not just co-ordinates

compute  
embedding  
of node  
from  
neighbours  
and  
store  
~~in~~ ~~in~~  
vector



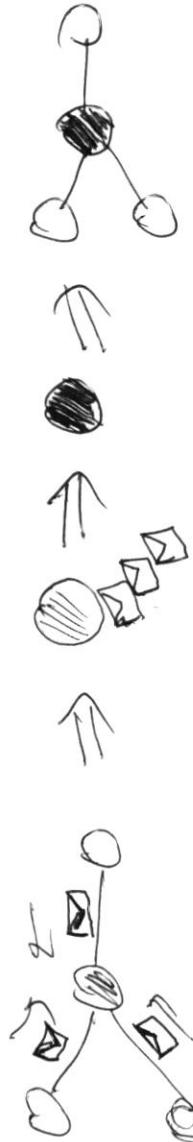
→ Neighbourhood-aware  
embeddings

e.g. if two nodes share many neighbouring nodes  
these two nodes should be close to each other in co-ordinates

Nodes: people

Node features: age, net worth  
edges: in phone contacts  
edge features: no. of phone calls in last year

② 2. Neural Message Passing framework



②  
2.1.2



~~This new feature is in addition to features 'eye' and 'network' of the individual.~~

After one/more iterations of neural message passing, the updated node has not just the 'eye' and 'network' features but also a new feature ~~that~~ that is computed by aggregating the features of the ~~messages~~ <sup>features of the messages</sup> from neighbours and then updating the previous value of this additional feature.

Newer Message Passing has 3 main functions:-

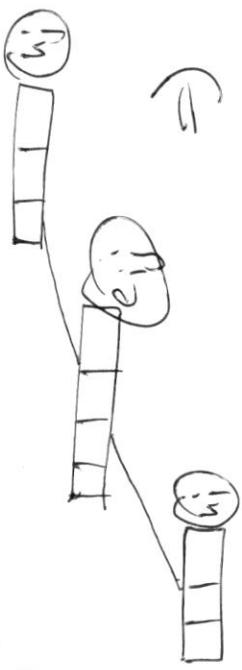
- ① Message function: → Computes messages from using node/edge features
- ② Aggregation fu.: → Combines all the ~~messages~~ from neighbouring nodes in a fixed-length vector,
- ③ Update fu.: → Computes the new value of embedding using aggregated messages and the old node embedding.

Message fu.

$$m_{ij}^{(k)} = \text{ML}(h_i^{(k)}, h_j^{(k)}, e_{ij})$$

→ message from node 'j' to 'i' in k<sup>th</sup> iteration of message passing

→ Message fu.



VI. 3

Example of message  $f_\alpha$  from literature

$$w_{ij}^{(k)} = h_j^{(k)} \rightarrow \text{Copy of the neighbor's features}$$

$$\rightarrow w_{ij}^{(k)} = \frac{1}{c_{ij}} h_j^{(k)} \leftarrow \text{Normalized (scaled by constant)}$$

$$\rightarrow w_{ij}^{(k)} = \frac{1}{c_{ij}} h_j^{(k)} \leftarrow \text{feature of neighbor.}$$

$$\rightarrow w_{ij}^{(k)} = \alpha \left( h_i^{(k)}, h_j^{(k)} \right) \leftarrow \begin{array}{l} \text{Scale neighbors' } \\ \text{features by a } \\ \text{constant} \end{array}$$

Attention  $f_\alpha$

The neighbors of node  $i$  are  $\{j_1, j_2, \dots, j_n\}$ .  
 Based on the similarity between the features of node  $i$  and  $j$ , the attention score  $s_{ij}$  is calculated.

The attention score  $s_{ij}$  is scaled by  $\alpha$ .

Then the feature  $h_j$  is combined with the feature  $h_i$ .

$s_{ij} = \text{score}(h_i, h_j)$  is scaled by  $\alpha$ .

VI.4

$$\text{Aggregation fn. } m_i^{(k)} = \bigoplus_{j \in N_i} m_{ij}^{(k)}$$

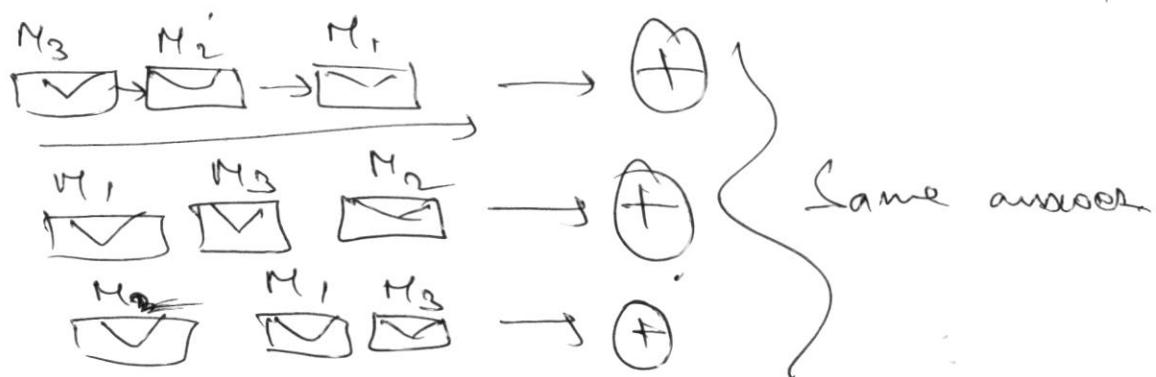
$\rightarrow \bigoplus$  aggregates all the messages  $m_{ij}^{(k)}$  coming from ~~one~~ the  $N_i$  neighbours of  $i^{\text{th}}$  node.

### Properties of Aggregation fn.

~~(1) Permutation invariant~~ gives the same answer regard

Two fn. should give the same answer regardless of:

- (1) No. of neighbouring nodes  $\leftarrow$  Fixed length representation
- (2) Order of the i/p messages  $\leftarrow$  Permutation invariant.



### Aggregation Function Examples:

Fixed length rep.  
Permutation invariant

$$(1) m_i^{(k)} = \sum_{j \in N_i} m_{ij}^{(k)} \quad (\text{Sum})$$

$$(2) m_i^{(k)} = \frac{1}{|N_i|} \sum_{j \in N_i} m_{ij}^{(k)} \quad (\text{Average})$$

$$m_i^{(k)} = \max_{j \in N_i} (m_{ij}^{(k)}) \quad (\text{Max})$$

v1.5

~~Different kinds~~Update fn.

$$h_i^{(k+1)} = \phi(h_i^{(k)}, m_i)$$

current value  
of embedding

aggregated  
message

Examples of Update fn.

$$\textcircled{1} \quad h_i^{(k+1)} = g(W^{(k+1)} m_i^{(k)})$$

$$\textcircled{2} \quad h_i^{(k+1)} = \phi(W_{\text{self}}^{(k+1)} h_i^{(k)} + W_{\text{neigh}} m_i^{(k)} + b^{(k+1)})$$



$$h_i^{(k+1)} = g(W^{(k+1)} \text{Concat}(h_i^{(k)}, m_i^{(k)}))$$

aggregated message

Unlike the 1st  
Update fn.

this fn. doesn't  
just aggregate the  
feature of the node  $h_i$

with the incoming  
message

Instead  $h_i$  treats the  
message as separable

Feature  $h_i$  and concatenated it  
with the aggregated  
message

Message Passing Framework:-

$$h_i^{(k+1)} = g\left(W^{(k+1)} \sum_{j \in N_i} \frac{1}{c_{ij}} h_j^{(k)}\right)$$

message fn.

aggregat fn.

Update fn.

VI.6

GraphSage:

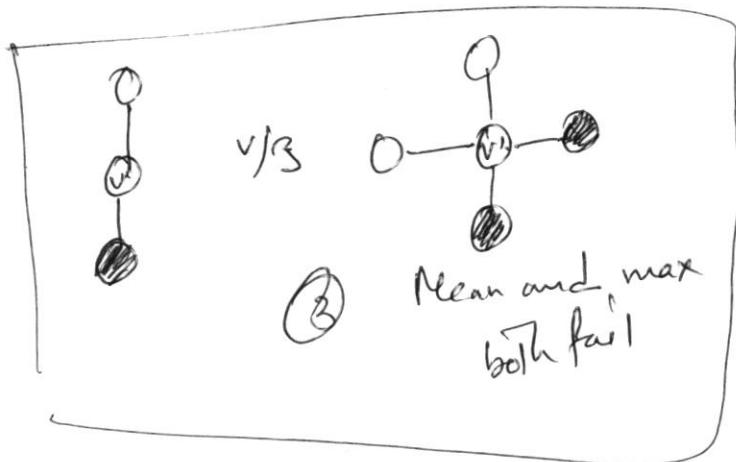
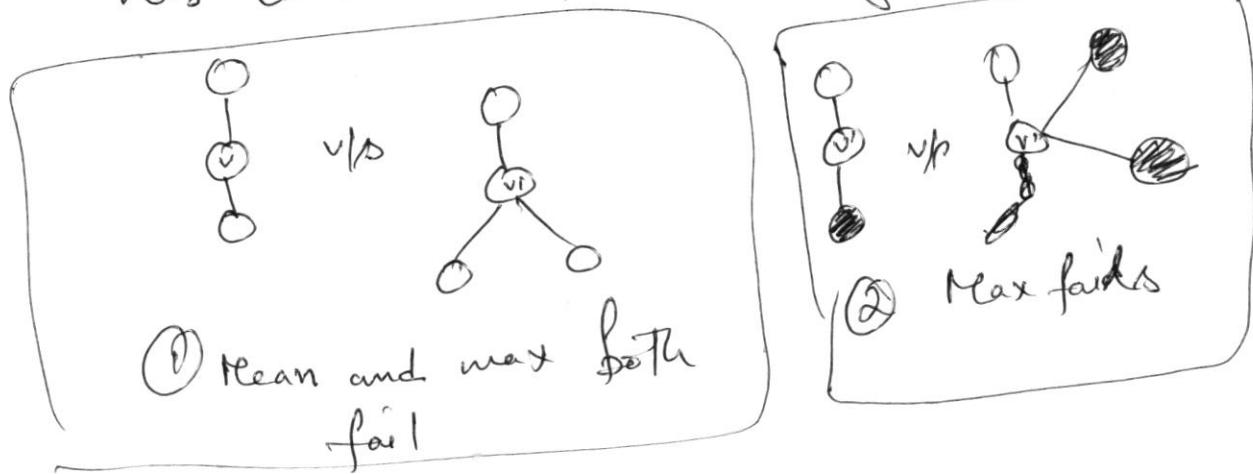
Function to  
compute  
neighbourhood-aware  
node embeddings

$$h_i^{(k+1)} \Rightarrow \left( W^{(k+1)} \text{ CONCAT } \left( h_i^{(k)}, \frac{1}{|N_i|} \sum_{j \in N_i} \right) \right)$$

→ Different choices for Message, Aggregate and Update  
functions are applicable depending on the application / data

Limitations: —

- ① Repeating the message passing ~~from~~ <sup>too many</sup>  
multiple times will cause the well-known "over-smoothing" problem, which makes the node embeddings become a self-similar blob.
- ② Aggregate fn. applied to pair-wise messages has limited ability to distinguish certain structures



**Permutation Symmetries**

7  
July

## Symmetries

~~GOAL: FIND A FUNCTION~~

WAT C. inherent order

1) Nodes from a graph

we're  
watching out for you

THB  
Func  
- Not

→ 2) Creating an adjective  
an adverb

the nodes in the graph

Hust No.  
DEPEND ON OUR  
RARY ORDER IN  
The import

ansigens offen

... depends on the

or input

→ 3) Traditional ML will  
i.e. if you shuffle the columns  
it won't work.

features, the model will give the same answer.

We want a model that  
regards the ordering of  
orderings as arbitrary.

$\lambda^2$  in the  
adjacency  
matrix

~~A simple graph will model the learnt different predictions for different orderings of features.~~

do  
of the  
the  
name  
and

by one

essentially

nodes give result

Different ordering (arbitrary)

Office in different adjacency matrix with different entries.

series of nodes

→ An MLO

trained with  
advice

will be  
a lesson

v2.2

Sol<sup>n</sup>: to the above problem

- ① Build these permutation symmetries into your training graphs.
- ② Train the MLP on all possible orderings/perm. of the nodes of all the training graphs.
  - ↳ Very expensive
    - ① So many comb's. poss'
    - ② So many training samples.

This problem of symmetry also exists in CNNs.  
e.g. A face in top left corner of an image is same as a face in top right/bottom left corner of the image.  
→ Can't train the net by providing all possible poss'ns. of possible poss'ns. of the image.  
overcame the problem i.e. scanning a filter across the image by convolution  
This scanning property puts symmetry into the model itself.

B) Permutation Invariance: A node's neighbourhood is a set of ~~fixed~~ nodes. The input ordering, but get the same result every time.

(Order information is ignored)

Whatever be the internal order of the elements of the set, the o/p remain the same (impts)

$f(Px) = f(x)$

It essentially maps a set of ips to an o/p. It doesn't depend on the ordering.

Permutation of the order of input

Permutation Equivariance: - if you shuffle the i/p

(Order information is maintained) You have to shuffle the o/p in the same way to get the same result.

$$f(Px) = Pf(x)$$

RBFB  
Raff V2.6

Example: Three

people: James, Sally and Mary.



James is socially connected

Soln: ①

for permutation symmetry problem

Permutation Invariance Using Aggregation f.

Aggregation like SUM | AVG | MAX | MIX

\* \* \* These aggregation functions are examples of permutation invariant functions

$$\sum (S - J - M) = \sum ((J - M) + S) \\ = \sum ((J - S) + M)$$

Similarly, for AVG / MAX / MIN functions

v2.4

Soln:- ② Permutation Equivariance Using Node  
application of a fn.

$$f(Px) = Pf(x)$$

Permutation Equivariance maintains the ordering information such that - it doesn't matter whether you permute before or after.

→ One way to do/achieve this is by:  
Applying the same fn. to every node individually  
 $n_i = \phi(n_i)$

$$f(Px) = Pf(x)$$

$$\Rightarrow \phi(P_{(3,1,2)}x) = \begin{bmatrix} -\phi(n_3) - \\ -\phi(n_1) - \\ -\phi(n_2) - \end{bmatrix} = P_{(3,1,2)}\phi$$

Order of i/p node  $\equiv$  Order of o/p nodes

~~④~~  $f(S-S-M) = S-\boxed{S}-M$

~~④~~  $f(S-M-S) = \boxed{S}-M-S$

~~General~~

i/p order  $=$  o/p order

fn. is applied on each input. whatever be the internal ordering of the elements of the set, the o/p is the

~~Graph is a collection of nodes on~~

Graph → Set of Notes + Edges

Set of  
Nodes

+ Bdgcs

of Nodes  $\equiv$  Set }

Collection of Nodes  $\equiv$  Set  
 Collection of Edges  $\equiv$  Adjacency Matrix

Edges  $\equiv$  Adjacency Matrix

Permuting nodes = Permuting elements of a set  
 Permuting edges = rows and columns of the adjacency matrix

Updated Defn.

# Permutation Invariants

Because not just a set of adjacencies in a matrix form and Permutation Equivalence

A graph with 16 nodes has  $D_{max}$  permutation invariance

Permutation Equivalence

Adjacency matrix

$$\begin{bmatrix} S & J & N \\ \hline S & 0 & 1 & 0 \\ J & 1 & 0 & 1 \\ N & 0 & 1 & 0 \end{bmatrix}$$

Adjacency Matrix

$$\begin{bmatrix} J & S & N \\ \hline J & 0 & 1 & 1 \\ S & 1 & 0 & 0 \\ N & 1 & 0 & 0 \end{bmatrix}$$

$f(Px, PAP^T) = Pf(X)$

Graph:

```

graph TD
    Suman --> James
    James --> Suman
    James --> Nick
    Nick --> James
  
```

Matrices:

- T:  $\begin{bmatrix} 0 & 0 & 1 \\ 0 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix}$
- J:  $\begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \end{bmatrix}$
- S:  $\begin{bmatrix} 0 & 0 & 1 \\ 0 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix}$
- N:  $\begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}$

V2.6

~~V2.1~~ Different matrices for adjacency matrices represent the same graph.

Any row and any column in an adjacency matrix can be assigned to any node of the graph.

V2.3

Permutation Invariance:

$$f(\{x_j : j \in N_i\}) = h_{N_i}$$

Set of nodes  
in a neighbourhood

embedding that  
represents  
the entire  
neighbourhood

→ So, i/p to a permutation invariant fn. is a set.

Permutation Equivariance:

$$h_i = f(x_i)$$

individual node

Node-level op.

In page ~~V2.4~~ we had:

$$\phi(X) = \begin{bmatrix} -\phi(x_1)- \\ -\phi(x_2)- \\ -\phi(x_3)- \end{bmatrix}$$

Function to create node embeddings

But, a graph is not just a set of nodes, but also has edges. So, the function has to take the edges too as i/p.

EXAMPLES:-  
GRAPH SAGG in V1.6

Ans

Fn. to compute  
"neighbourhood" aware  
node-embeddings

Adjacency  
matrix

permutation equivariance function

$$\phi(X, A) = \begin{bmatrix} -\phi(x_1, \{x_j : j \in N_1\})- \\ -\phi(x_2, \{x_j : j \in N_2\})- \\ -\phi(x_3, \{x_j : j \in N_3\})- \end{bmatrix}$$

$$\phi(x, A) = \begin{bmatrix} -\phi(x_1, \{x_j : j \in N_1\}) - \\ -\phi(x_2, \{x_j : j \in N_2\}) - \\ \vdots \\ -\phi(x_n, \{x_j : j \in N_n\}) - \end{bmatrix}$$

~~GNNs that computes neighbourhood-aware node embedding without varying with~~

~~Message for Update fn. (Permutation equivariant fn.)~~

SUMMARY:- GNNs are functions that produce node-level outputs (taking neighbourhoods into account) that do not depend on our arbitrary ordering (of nodes in adjacency matrix)

i.e. permutation equivariant fn.  $\phi$

$$\phi(x, A) = \begin{bmatrix} -\phi(x_1, \bigoplus_{j \in N_1} x_j) - \\ -\phi(x_2, \bigoplus_{j \in N_2} x_j) - \\ \vdots \\ -\phi(x_n, \bigoplus_{j \in N_n} x_j) - \end{bmatrix}$$

~~Only when the i/p to  $\phi$  is constant/invariant, we can expect an overall equivariant  $\phi$ .~~

$\bigoplus$   $\leftarrow$  must be permutation invariant. Otherwise, the fn.  $\phi'$  won't produce equivariant o/p.

Each ' $\phi$ ' on ~~the~~ input node and its neighbourhood ~~or~~ could produce different o/p depending on the ordering of neighbouring node-messages.

Q2.8

So, we want a fn. (a GNN) :-

- ① that doesn't depend on the ordering of neighbouring nodes
- ② But does ~~not~~ maintain the ordering of nodes. Otherwise we don't know which ~~node~~ ('embedding' is computed for which 'node' of the graph).

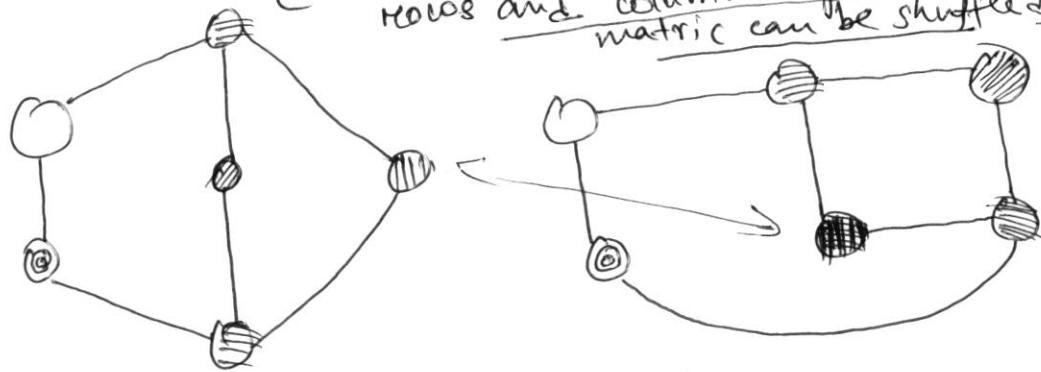
Soln. → Get node-level embeddings that respect node ordering (Permutation Equivariant) by applying a fn. to each node neighbourhood (Permutation Invariant w.r.t to nodes of neighbourhood), i.e. doesn't depend on the order of the nodes in the neighbourhood set.

V3.1

## Graph Isomorphism

Two graphs are isomorphic if they are essentially identical but perhaps have been assigned a different ordering (i.e. different adjacency matrices)

→ Graphs don't have an inherent order and any order that we impose is arbitrary (e.g. by putting nodes into an adjacency matrix) [There are many ways in which the rows and columns of an adjacency matrix can be shuffled]



Isomorphic graphs

Q: Does Graph Isomorphism matter for GNNs? Why?

→ If  $G_1$  and  $G_2$  are isomorphic graphs, does your GNN give the same o/p?

→ If  $G_1$  and  $G_2$  are not isomorphic, does your GNN give different o/p's?

~~Q:~~  $G_1(x_1, A_1)$  and  $G_2(x_2, A_2)$  ~~are~~ <sup>should be</sup> isomorphic if

there's a  $P$  such that:-

$$P^T A_1 P = A_2 \text{ and } P x_1 = x_2$$

i.e. if you can re-order one to obtain the other

~~How?~~

v3.2 How to find if two graphs  $G_1$  and  $G_2$  are isomorphic

Soln: ① Brute-force: Test every possible permutation of an adjacency matrix of  $G_1$  and see if until you obtain the adjacency matrix of  $G_2$  or you just run out of permutations.

But the no. of permutations scales with the factorial of the no. of nodes

$$\text{e.g. } 100! = 9 \times 10^{157}$$

So, Brute-force is not a practical approach.

Soln: ② Weisfeiler - Lehman(WL): A family of algorithms for testing whether two graphs are isomorphic.

1-WL steps:-

For each node in the graph,

Step ① Assign an initial label (say, 0) to each node of the graph

Step ② Assign a new label by mapping each tuple

(current label, set of neighbouring labels) to a new label

$$\Rightarrow \ell^{(i+1)}(u) = \left\langle \left( \ell^{(i)}(u), \{ \ell^{(i)}(v) : v \in N(u) \} \right) \right\rangle$$

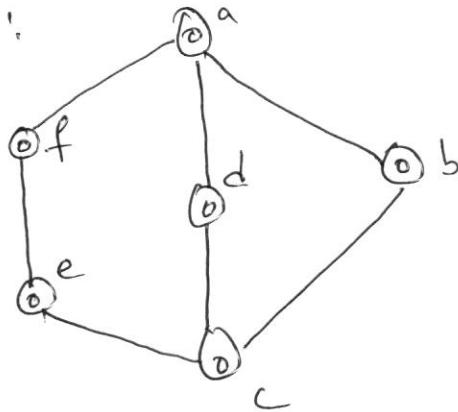
Step ③: Partition the graph into node subsets that share the same colour

Step ④: Repeat steps 2-3 until the graph partition stop changing

Perform the

Example:

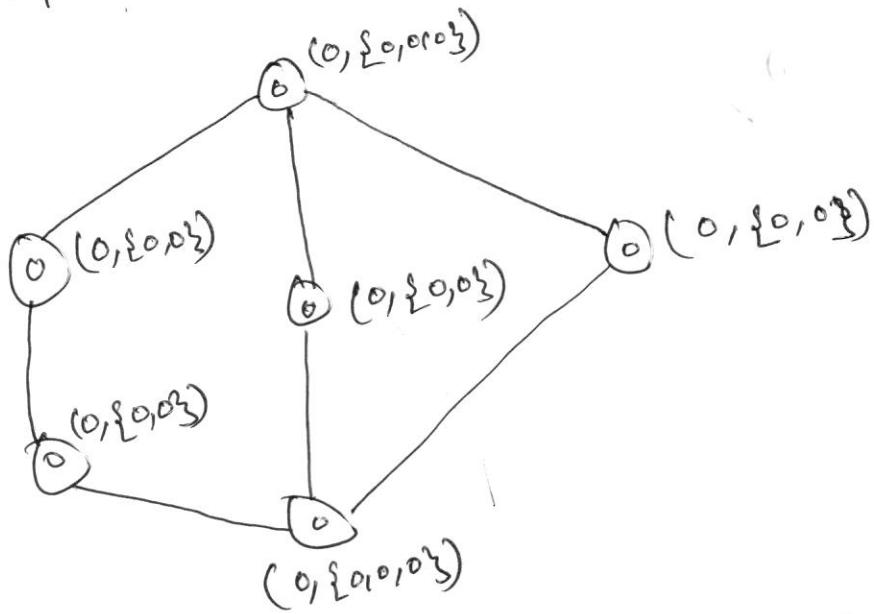
Step 1:



All nodes have same label '0'.  
Label '0' is represented by colour  
blue

Step 2: Create a tuple for each node

Tuple: (current label, set of neighbouring labels)



Step 3.

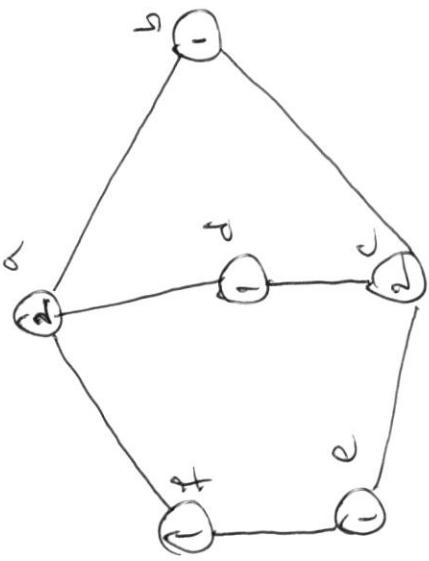
Assign a ~~unique colour~~ to each node  
Initial —  blue  
 $(0, \{0,1,0\})$  —  yellow  
 $(0, \{0,1,0,0\})$  —  blue

Partition the graph based on node labels/colours

13.1 Partition  
Recurrence

The graph based on three new labels / colours

Step 1:



Partition

[2]

[1]

{a, b, c}

{d, e, f}

Repeat the above steps on the newly labeled graph

Step 2:

(2, {1, 1, 1, 1})

(1, {1, 1})

(1, {2, 2})

(1, {2, 2, 2})

(1, {2, 2, 2})

(2, {1, 1, 1})

(2, {1, 1, 1})

(1, {1, 1})

(1, {1, 1})

(1, {1, 1})

Unique labels, colour

black

green

red

[3]

[4]

[5]

Unique Tuples,

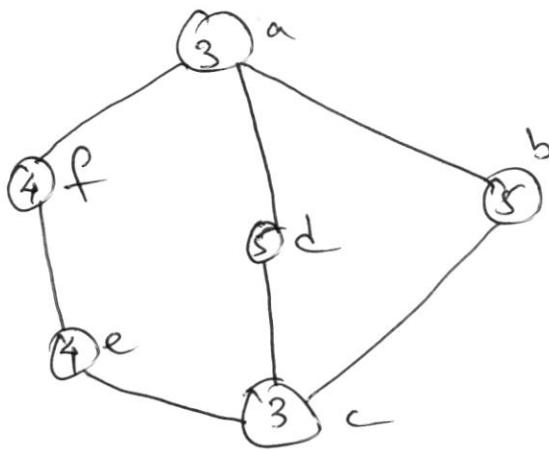
(2, {1, 1, 1})

(1, {2, 1})

(1, {2, 1, 2})

(1, {2, 1, 2, 3})

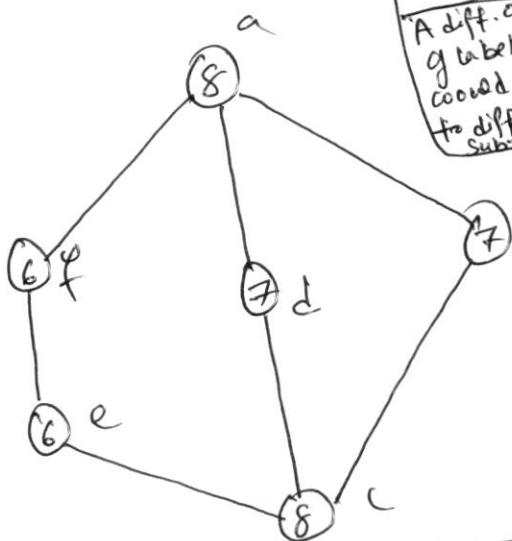
Step (5): Partition the graph based on the new labels.



Partitions:

<u>3</u>	<u>4</u>	<u>5</u>
{a,c}	{f,e}	{b,d}

Step (6) / (7):



Note: The elements of these sub-sets could vary depending on how the nodes are assigned labels.  
A diff. color labelling could lead to diff. sub-sets.

Partitions:

So, you can see the two partitions are identical.  
The two partitions lead to the same sub-sets of nodes.  
So, the algorithm has terminated/converged.

If two graphs or more graphs  $G_1, G_2, G_3, \dots$  converge to the same sub-sets of nodes i.e. partitions, it implies they are all isomorphic.

v3.6

If ~~three~~ two or more graphs  $G_1, G_2, G_3, \dots$  converge with the ~~same~~ in the LWL test, with the same no. of sub-sets - each containing the same no. of nodes - , they are isomorphic.

Two partitions for, while comparing the partitions of two or ~~if~~ more different graphs, the labels of the elements (nodes) in their subsets don't matter

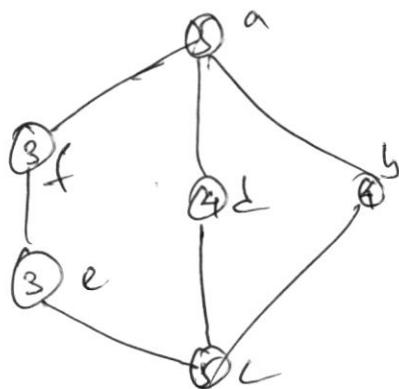
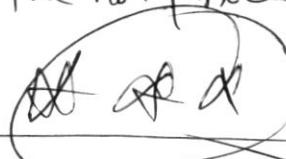
~~If~~ The labels will vary because isomorphic graphs differ in the ordering of the node labels.

~~Only the structural~~

in the partition

The exact labels don't matter.

Only the structural information - The no. of nodes per sub-set (partition) - matters



Partitions

[3] [4] [5]  
 $\{a, b\}$   $\{c, d\}$   $\{e, f\}$

[3] [4] [5]  
 $\{a, b\}$   $\{e, f\}$   $\{c, d\}$

Some graphs which are not isomorphic will still pass the 1-WL test.

→ Failing the 1-WL test guarantees that the two graphs are not isomorphic

→ Passing the 1-WL test doesn't guarantee that the two graphs are isomorphic. They may still be non-isomorphic.

The following eqn. is in Step 2 of the 1-WL test:

$$\ell^{(i+1)}(u) = \mathcal{L} \left( (\ell^{(i)}(u), \underbrace{\{\ell^{(i)}(v) : v \in N(u)\}}_{\text{new label}}) \right)$$

(Putting new labels to nodes based on their current label and the labels of the neighbouring nodes)

And, the GIN fn. is on page v2.7:-

$$\phi(x, A) = \begin{bmatrix} -\phi(x_1, \bigoplus_{j \in N_1} x_j) - \\ -\phi(x_2, \bigoplus_{j \in N_2} x_j) - \\ \vdots \\ -\phi(x_n, \bigoplus_{j \in N_n} x_j) - \end{bmatrix}$$

These two eqn. fns. bear similarity in the sense

- that - <sup>both</sup> :-
- ① Are applied on each node of the graph
  - ② Take into ac/c the neighbouring nodes
  - ③ Take into ac/c the current node

v3.8

$$\text{1-WL} := \ell^{(i+1)}(u) = \sigma \left( \sum_{v \in N(u)} \ell^{(i)}(v) \right)$$

So,

$$\text{GRAPHSAU}^0 - h_i^{(k+1)} = \sigma \left( W^{(k+1)} \text{CONCAT} \left( h_i^{(k)}, \frac{1}{|N_i|} \sum_{j \in N_i} h_j^{(k)} \right) \right)$$

→ 1-WL is effectively doing message passing iterations.

→ 1-WL first assigns a new label to each unique tuple i.e.

1-WL first assigns a new ~~label~~ label to each kind of (node + neighbourhood) multi-set

But, MP-GNN (message passing) can be less discriminative than such "deliberately mapping of every unique multi-set

e.g. Average  $\{\{1, 1, 1\}\} = \text{Average } \{\{1, 1\}\}$ . [Average as aggregate fn.]  
↳ Different neighbourhoods but same o/p.

But, 1-WL would assign a different label to  $\{\{1, 1, 1\}\}$  and  $\{\{1, 1\}\}$ .

So, the discriminative power of MP-GNN  $\leq$  Discriminative power of 1-WL

1-WL provides an upper-bound for the discriminative power of MP-GNN.

Two graphs may not be ~~isomorphic~~ isomorphic. But, there can be a measure of similarity b/w the two graphs.

How to measure this similarity?

- D ① Perform the 1-WL test on a ~~each~~ graph
- ② Create a list of each type of labels encountered along with the no. of nodes having ~~at~~ those labels until the convergence of the iterations of 1-WL test.

e.g.

	<u>Label maps</u>	<u>No. of nodes</u>
for initial	$\rightarrow [0] \hat{x}_0$	6
$(0, \{0, 0\}) \rightarrow [1] \hat{x}_1$	4	
$(0, \{0, 0, 0\}) \rightarrow [2] \hat{x}_2$	2	
$(1, \{1, 2\}) \rightarrow [3] \hat{x}_3$	2	
$(1, \{2, 2\}) \rightarrow [4] \hat{x}_4$	2	
$(2, \{1, 1, 1\}) \rightarrow [5] \hat{x}_5$	2	

- ③ Create a vector of these features:
- $$6\hat{x}_0 + 4\hat{x}_1 + 2\hat{x}_2 + 2\hat{x}_3 + 2\hat{x}_4 + 2\hat{x}_5$$

- ④ Perform these steps for the other graph and create its vector.
- ⑤ Take a dot product of ~~these~~ these 2 vectors. The result is a measure of similarity.

41

How Powerful are GNNs?

④ 1-WL test is and GNN are analogous.

 The message passing step collects all the neighbours' labels into a multi-set (i.e. a set where the same item can occur more than once)

~~in essence~~,

SUMMARY:- In ~~case~~ In 1-WL test, each unique neighbour hood gets a label and additional ~~iteration~~ considers larger neighbour hoods.  
Similar to The iterations of the message passing.

Each successive ~~iteration~~ message passing iteration of GNN  
~~includes~~ produces "neighbourhood-aware" embeddings  
that consider a larger neighbourhood.

So, the question "How Powerful are GNNs?" can be reduced to The question "Can a GNN map each ~~diff~~ unique neighbourhood to a different embedding (label)?"

§ 4 REMINDERS — A GNN must ~~not~~ produce  
“neighbourhood-aware” embeddings for each  
node of the n/w without being dependent ~~of~~ on the  
order of the nodes in the ~~a~~ neighbourhood of  
~~a node~~ the current node.

V4.2

$G_1$  and  $G_2$  are two distinct graphs.

Let's say, ① GNN 'A' distinguishes ~~is~~  $G_1$  and  $G_2$ . But 1-WL test

↓

② 1-WL produces identical partitions / sub-sets of no.

↓

③ 1-WL ~~obtained~~ ~~the same~~ identically labeled ~~sets~~  
graphs for  $G_1$  and  $G_2$ .

↓

④ GNN 'A' will fundamentally be given these identical  
label multi-sets and its Aggregate and update fn. will  
map those to identical o/p and therefore can't  
distinguish them too.

So, if 1-WL test fails, GNN will fail to

→ 1-WL is the upper bound for GNN

When Do GNNs meet this Upper Bound ~~set by 1-WL test~~

Ans. → { When the (Aggregate) and (Update) fn.s are  
Injective functions.

(Injective fn. → Maps distinct i/p s to distinct o/p s)

GIN  
(Graph Isomorphism Network)

Distinct multisets of node features are mapped to distinct embeddings.

An Injective Aggregate fn.

- 1) Do one-Hot-Bencoding for each ~~the~~ unique label in The graph
- 2) Represent each label as One-hot-encoded vector.
- 3) Add all Add(Sum) all ~~the~~ the vectors (corresponding to each label) in a multi-set set.

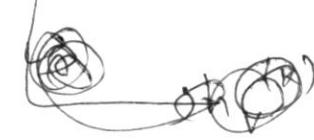
12)  $L = U \Lambda U^T$ ,  $\Lambda = \text{diag}([\lambda_0, \dots, \lambda_{N-1}])$ ,  $U = \text{matrix of eigen-vectors as columns}$

$\{a, b\}$  $\boxed{1 \ 1}$  $\{a, a, b, b\}$  $\boxed{2 \ 1 \ 2}$ 

$\therefore$  The above two neighbours can be distinguished by this injective fn.

On the other hand, SUM and AVERAGE fns  
can't distinguish these two neighbourhoods.

Next, we need an injective Update fn.



$$h_v^{(k)} = \phi((1 + \epsilon^{(k)}) \cdot h_v^{(k-1)} + \sum_{u \in N_v} h_u^{(k-1)})$$

node's previous representation

Why do we need the constant ' $\epsilon$ '?

example:  $\phi(h_1, \{h_2\})$  and  $\phi(h_2, \{h_1\})$

$$\underline{h_1 + \epsilon h_1 + h_2} \neq \underline{h_1 + h_2 + \epsilon h_2}$$

If there were no  $\epsilon$ , or,  $\epsilon=0$ ,

$$\begin{cases} \text{LHS} = h_1 + h_2 \\ \text{RHS} = h_1 + h_2 \end{cases}$$

So,  ~~$\phi(h_1, \{h_2\})$  and  $\phi(h_2, \{h_1\})$~~  could become ~~indistinguishable~~.

GIN :- (Graph Isomorphism Network)

$$h_v^{(k)} = \text{MLP}^{(k)} \left( (1 + \epsilon^{(k)}) \cdot h_v^{(k-1)} + \sum_{u \in N(v)} h_u^{(k-1)} \right)$$

[one-hot-encoded]

~~1~~~~V4.4~~

GIN pr

V5.1

The Graph Laplacian

→ Degree matrix of a graph is ~~the~~ a diagonal matrix in which diagonal element represents the degree of a node of the graph.

↳ which diagonal element represents which node is arbitrarily chosen

→ Degree of a node is the no. of edges attached to ~~one~~ a node

→ In an undirected graph, each loop increases the degree of a vertex by 2.

~~1~~ ~~2~~ ~~3~~  
Degree of node  
 $\Rightarrow$  Sum of all ~~edges~~  
Mows along ~~row~~  
~~of~~ ~~column~~ (node)  
of the Adjacency matrix.

→ In a directed graph, the term 'degree' may refer to either to 'indegree' (no. of incoming edges at each vertex) or 'outdegree' (no. of outgoing edges at each vertex).

→ Laplacian ( $L$ ) =  $D - A$ , where  $D \rightarrow$  Degree matrix } of size  $n \times n$   
 $A \rightarrow$  Adjacency matrix } of size  $n \times n$

$$L_{ij} = \begin{cases} d_i & \text{if } i=j \\ -1 & \text{if } (i,j) \in E = \delta_{ij} d_i - A_{ij} \\ 0 & \text{otherwise} \end{cases}$$

→ The diagonal elements of any adjacency matrix are zero. Rest of the elements can be 1s or 0s depending on whether there is ~~an~~ an edge b/w the nodes represented by the corresponding ~~row, column~~ pair.

→ Graph Laplacian is a pre-requisite for understanding graph convolutions.

$$(L_n)_{ij} = \sum_j (\delta_{ij} d_i - A_{ij}) \otimes x_j.$$

(Laplacian operating on a graph node)

( $j$ th element of this product)

$$\stackrel{\text{for } i=j}{=} d_i - \sum_j A_{ij} x_j$$

1. Laplacian ( $L$ )

→ where there are connections

else

Adjacency matrix ( $A$ )  
1 for connections  
0 for others

3.  $D \rightarrow$   
degree for diagonal elements

12)  $L = U \Lambda U^T$ ,  $\Lambda = \text{diag}([x_0, \dots, x_{N-1}])$ ,  $U = \text{matrix of eigen-vectors as columns}$

N5.2

Degree of  $i$ th node of a graph = Sum of all the elements in the  $i$ th column of the adjacency matrix ( $A$ ) of the graph.

$$\Rightarrow d_i = \sum_j A_{ij}$$

why  $\sum$  came from  $\mathbb{Z}$   
came from  $L_2$

$\sum_j$  means across all "j's" for coming

$i$ th node/element  $(L_n)_i = \sum_{jk} (\delta_{ij}d_j - A_{ij})x_j$

$L_n \rightarrow$  Laplacian operator applied on all nodes of a graph

$A_{ij} \rightarrow$  zeros out all the elements except the connections/edges between nodes  $(i, j)$

$\therefore \sum_i$  over all nodes can be replaced by  $\sum_j$  over just the connections

$$= d_i x_i - \sum_j A_{ij} x_j$$

$$= \sum_j A_{ij} (x_i - x_j)$$

$$= \sum_{(i,j) \in E} (x_i - x_j)$$

$E \leftarrow$  Set of all edges in the graph

Difference b/w a node's value and its neighbour's values

$$(L_n)_i = \sum_i (\delta_{ii}d_i - A_{ii})x_i$$

$$= d_i x_i$$

'L' is discrete equivalent of ' $\nabla$ '.

In physics, Laplacian operator  $\nabla^2$  is "average difference b/w a point and a small sphere around that point"

In discrete setting, operating on graph nodes,

$\int_0^1$  is replaced by

V5.3

$$(L_x)_i = \sum_j A_{ij} (x_i - x_j)$$

←  $A_{ij} > 0$  is everything than the connection  
So,  $\sum_j A_{ij} (x_i - x_j)$

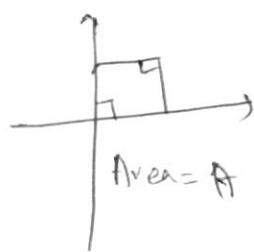
$$= \sum_{(i,j) \in E} (x_i - x_j)$$

Sum of the Differences b/w the embeddings of the node we care about - ( $i^{\text{th}}$  node) and all the embeddings of all its neighbouring nodes ( ~~$j^{\text{th}}$  node~~) (denoted by  $x_j$ ).

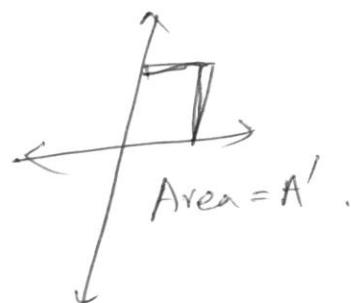
Basic linear

⊗⊗ Refresher on Basic linear Algebra

Determinant →



$\rightarrow$   $A$   
linear transform-  
ations by  
applying  
transforms on  
matrix  $A$ .



Determinant = Scaling factor of a chunk of area

⊗⊗⊗

of by applying a linear  
matrix transformation.

→ Determinant is a measure of the transformation brought about by applying a  $3 \times 3$  matrix  $A$ .

12)  $L = U \Lambda U^{-1}$ ,  $\Lambda = \text{diag}((\lambda_0, \dots, \lambda_{N-1}))$ ,  $U = \text{matrix of eigen-vectors as columns}$

v5.4

A determinant of 0 means the area/space was collapsed to a point/~~line~~ (a reduction in dimensionality) by applying the transformation matrix.

⇒ In such a case, there is ~~is no inverse~~  
transformation.

A negative determinant ~~means~~ means the space is

flipped (RHS → RHS, Right hand ~~thumb rule~~ → left hand ~~thumb rule~~)

( $\rightarrow$  coordinate system

~~thumb rule~~

coordinate

systems

~~Ax = Ax~~

Eigen Vectors tell you the axes ~~at~~ along which there is

just stretching and no rotation ⇒ Vectors that  
are mapped exactly back on themselves, but  
perhaps shorter or longer.

~~A~~ ~~or~~ ~~or~~

More the eigen value, greater the stretching along eigen vectors

$Ax = \lambda x$

$\lambda \rightarrow$  Eigen Value

$x \rightarrow$  Eigen Vector

More the squishing/squeezing along eigen-vector,  
less smaller the ~~Ax~~ eigen-value

The Eigen-value of an eigen-vector tells how much the eigen-vector is scaled

$\lambda$  and  $|A|$  (determinant) are analogous

$\lambda \rightarrow$  amount of scaling of eigen vector

$|A| \rightarrow$   $n \times n$  in a space/area

Product of all eigen values of a matrix 'A'

$= |A| =$  Overall amt. of scaling

~~A~~ ~~A~~ ~~A~~

v 5.5

$$Ax = \lambda x$$

① Transformation matrix 'A' is known.

$x \rightarrow$  Eigen vector

$\lambda \rightarrow$  Eigen Value

Find  $x$  and  $\lambda$ .

If a vector  $x$  is such that  $Ax = \lambda x$ , then  $x$  is an eigen vector if  $\lambda = \text{const}$

$$Ax = \lambda x$$

$$\Rightarrow (A - \lambda I)x = 0 \quad \text{--- (1)}$$

$$\Rightarrow x = 0 \quad \text{or} \quad A - \lambda I = 0 \quad \checkmark$$

$$x \Rightarrow |A - \lambda I| = 0$$

Find  $\lambda$ .

Plug ' $\lambda$ ' in (1) to get ' $x$ '.

Transformation matrix  $(A - \lambda I)$

Every row of a Laplacian sums to 'zero'.

~~Proof~~ Proof:-  $(L)_i = \sum_j L_{ij} \leftarrow \sum_j \text{ for constant } i \rightleftharpoons \sum_j \text{ occurs upon multiplying } L \text{ by a vector of } 1 = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$

$\nearrow \text{across all } j's \quad \nearrow \text{for constant } i$

$$= \sum_j (\delta_{ij} d_i - A_{ij})$$

$$= (\delta_{ij} d_i - A_{ij})_{i=j} + \sum_{i \neq j} (\delta_{ij} d_i - A_{ij})$$

'L' is a transformation matrix here

$$= (1 \cdot d_i - 0)_{i=j} + \sum_{i \neq j} (0 \cdot d_i - A_{ij})$$

\* \* \*

is a vector of  $1's$   
No. of  $1's$  depends on  
no. of nodes in graph

$|L| = 0$   
 $L$  has no inverse

$$d_i - \sum_{j \neq i} A_{ij}$$

$$= d_i - d_i = 0 \Rightarrow 1 \text{ is an eigenvector of } L \text{ with eigen-value } 0$$

$\Rightarrow 1$  is an eigenvector of  $L$  with eigen-value  $= 0$

\* \* \*

Sol:

1 is an eigen-vector of  $L$  with eigen-value  $= 0$ .

12)  $L = U \Lambda V^T$ ,  $\Lambda = \text{diag}([\lambda_0, \dots, \lambda_{N-1}])$ ,  $U = \text{matrix of eigen-vectors as column}$

55 v5.6

## Real Symmetric Matrix:

1) Real valued elements

2)  $i^{\text{th}}$  row =  $i^{\text{th}}$  column

Symmetric matrix

'L' is a  
real-symmetric  
matrix  $\star\star\star$

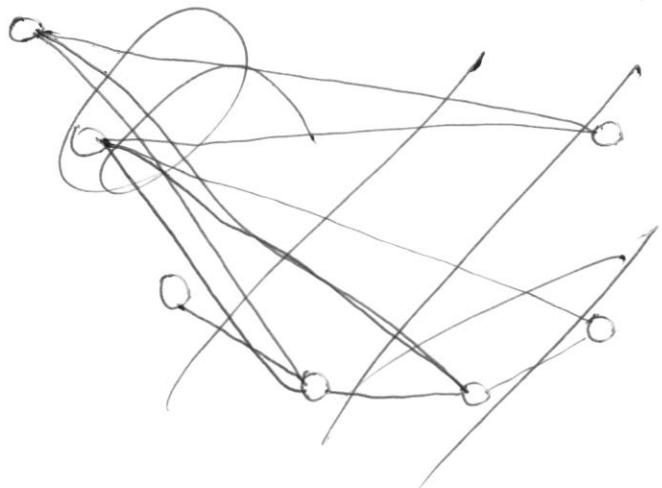
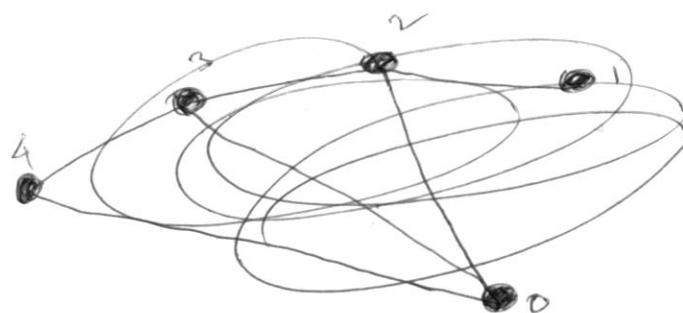
∴ Matrix is equal  
to its transpose

A Real Symmetric Matrix has real eigen-values  
and orthogonal eigen-vectors.  $\star\star\star$

All eigen-values of 'L' are non-negative  $\star\star\star$

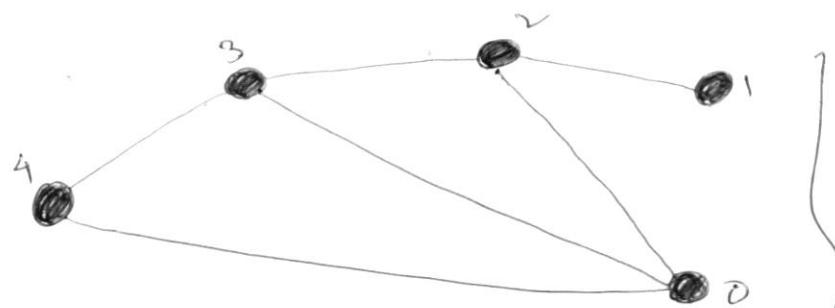
Then, how do we have negative  
 $\Rightarrow [L]$  can't be  $< 0$ . determinants.

'L' (Laplacian) is a Real Symmetric Matrix

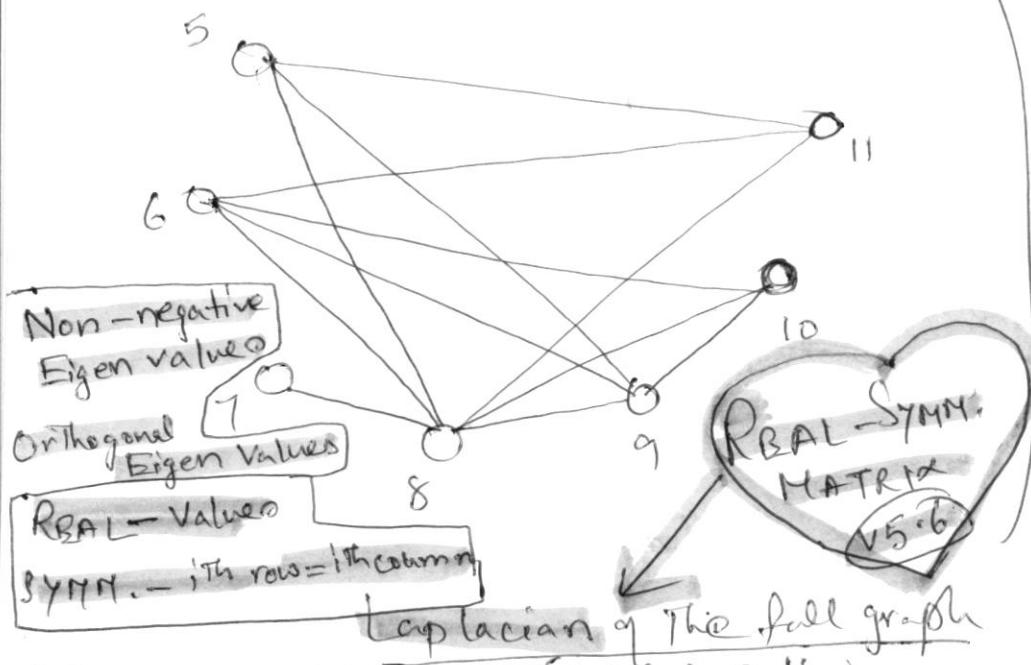


✓ 5.07

## Graph Partitioning



A graph with  
2 components



$d_0 - 3$	$d_9 - 4$
$d_1 - 1$	$d_{10} - 3$
$d_2 - 3$	$d_{11} - 3$
$d_3 - 3$	
$d_4 - 2$	
$d_5 - 3$	
$d_6 - 4$	
$d_7 - 1$	
$d_8 - 6$	

Check in,  
the following  
laplacian  
matrix

$$\sum d_{ii} = 0$$

	0	1	2	3	4	5	6	7	8	9	10	11
0	0	-1	-1	-1	-1	0	0	0	0	0	0	0
1	-1	0	0	0	0	0	0	0	0	0	0	0
2	-1	-1	0	0	0	0	0	0	0	0	0	0
3	-1	0	-1	0	0	0	0	0	0	0	0	0
4	-1	0	0	-1	0	0	0	0	0	0	0	0
5	0	0	0	0	0	0	0	0	0	0	0	$d_5$
6	0	0	0	0	0	0	0	0	0	0	$d_6$	
7	0	0	0	0	0	0	0	0	0	0	$d_7$	
8	0	0	0	0	0	0	0	0	0	0	$d_8$	
9	0	0	0	0	0	0	0	0	0	0	$d_9$	
10	0	0	0	0	0	0	0	0	0	0	$d_{10}$	
11	0	0	0	0	0	0	0	0	0	0	$d_{11}$	

The upper left half of the 'L' is the 'L' of the upper component of the graph.

The lower left half of the 'L' is fully black due to lack of any connection between the two components of the graph.

The upper right half of the 'L' is fully black (or zero) due to lack of any connection between the two components of the graph.

The lower right half of the 'L' matrix is the 'L' of the lower component of the graph.

Laplacian :-

- 1 for edge connections
- $\text{degree}(d_i)$  for diagonal elements
- 0 for others.

$$(2) L = U \Lambda U^T, \Lambda = \text{diag}((\lambda_0, \dots, \lambda_{N-1})), U = [u_0 | u_1 | \dots | u_{N-1}] \text{ as column}$$

5.8

## GRAPH PARTITIONING

~~We In the previous matrix of the previous page, we can treat it~~

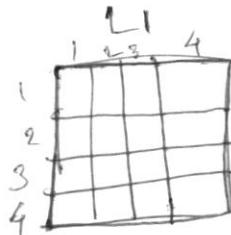
In the graph of the previous page, we can treat each of the two components as separate entities (graph) and create separate Laplacians for them. We can call them  $L_1$  and  $L_2$ .  $L_1$  forms the upper left half of the L matrix in the previous page.  $L_2$  forms its lower right half.

No. of eigen-vectors (of a Laplacian) with eigen-value = 0

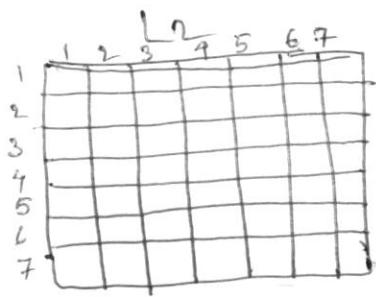


is the no. of disconnected components in a graph

If you don't know how many disconnected components there are in a graph, take the laplacian of the entire graph and figure out how many Eigen vectors have Eigen values = 0 and the value of those Eigen vectors will tell you about the nodes of each of the disconnected components.



$$L_1 \mathbf{1} = \lambda_1 \mathbf{1} \text{ where } \lambda_1 = 0, \text{ Here } \mathbf{1} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} \text{ a vector of 4 elements}$$



$$L_2 \mathbf{1} = \lambda_2 \mathbf{1} \text{ where } \lambda_2 = 0$$

Here  $\mathbf{1} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$  } A vector of 7 is

VS.9

GRAPH PARTITIONING

~~①~~  $X_{L1} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$  and  $X_{L2} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$

4 Is                                              7 Is'

were Eigen-vectors for  $L_1$  and  $L_2$  respectively

- ② Eigen vectors for the whole Laplacian  $L$  ~~and be~~  
can be derived from the above two Eigen-vectors.
- ||

By padding the above vectors of 1s with 0s. No. of 0s would be the remaining no. of nodes. For example, the 1st Eigen-vector for  $L$ , was 4 1s ~~as~~ while the total no. of nodes in the graph is 11. So, we will add 7 ( $= 11 - 4$ ) 0s. Similarly, to get ~~the~~ eigen-vector for  $L$  from the eigen-vector of  $L_2$ , we would add 4 ( $= 11 - 7$ ) 0s to it.

add

$$\textcircled{2} \quad X_{L1} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} \rightarrow X_L = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad \text{11 elements}$$

$$X_{L2} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} \rightarrow X_L = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix} \quad \text{11 elements}$$

## GRAPH PARTITIONING

In So, we ~~added~~ <sup>concatenate</sup> 0s at the end of  $x_{11}$  and at the beginning of  $x_{12}$  to obtain the Eigen vectors of  $L$ . ~~Why?~~

This can be generalized as: To create Eigen vectors for the whole Laplacian ' $L$ ' of a graph, we put 1s for nodes of one component and 0s for nodes of another component of the graph. We redraw the ' $L$ ' of the graph considered. This way we create multiple eigen-vectors - each one by putting 1s corresponding to nodes of one component and 0s corresponding to nodes of another component. To illustrate, we ~~create~~ redraw the  $L$  matrix for the previous graph.

	0	1	2	3	4	5	6	7	8	9	10	11	$x_1$	$x_2$
0	0	1	2	3	4	5	6	7	8	9	10	11	0	0
1	0	0	-1	-1	-1	0	0	0	0	0	0	0	1	0
2	0	0	1	0	2	0	0	0	0	0	0	0	1	0
3	0	-1	-1	0	-1	0	0	0	0	0	0	0	1	0
4	0	-1	0	-1	0	3	0	0	0	0	0	0	1	0
5	0	0	0	0	0	0	4	0	0	0	0	0	0	0
6	0	0	0	0	0	0	5	0	0	0	0	0	0	0
7	0	0	0	0	0	0	6	0	0	0	0	0	0	0
8	0	0	0	0	0	0	7	0	0	0	0	0	0	0
9	0	0	0	0	0	0	8	0	0	0	0	0	0	0
10	0	0	0	0	0	0	9	0	0	0	0	0	0	0
11	0	0	0	0	0	0	10	0	0	0	0	0	0	0

$$Lx_1 = \lambda_1 x_1 = 0 \quad (\because L_{11} = 0)$$

$$Lx_2 = \lambda_2 x_2 = 0 \quad (\because L_{22} = 0)$$

↑  
Is for  
nodes 0-4

↑  
Is for  
nodes 5-11

Eigen vectors created

These may have Eigen-Value

Eigen Vector created from the vector '1' for Laplacians of component graphs ~~of graph~~

V.S.U

## GRAPH PARTITION PARTITIONING

Even if the ordering of the nodes of the graph in V5.7 is shuffled, the L will be modified.



To create Eigen vectors with ~~all~~ Eigen values = 0, just follow rule 1s for nodes (as per new ordering) of one component of graph and 0s for the other components.

for example, the new Eigen vectors for the restuffed graph will look like:

$$\lambda_1 = \frac{1}{\sqrt{7}} \begin{bmatrix} 1 \\ 1 \\ 0 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix} \quad \lambda_2 = \frac{1}{\sqrt{5}} \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$\lambda_0 = \lambda_1 = 0$

To figure out the nodes of the disconnected components of a graph, i.e.

to figure out the partitions of a graph, i.e. :-

- ① Find the graph Laplacian by ordering nodes.
- ② Find out all the Eigen vectors of the graph.
- ③ Find out the Eigen vectors with eigen values of ~~set of~~.
- ④ The 1s in each of these Eigen vectors correspond to each disconnected component. The nodes of each component.



5.12

## Positive Semi-Definite

$$\begin{aligned}
 x^T L x &= \sum_{ij} L_{ij} x_i x_j \\
 &= \sum_{ij} (S_{ij} d_i - A_{ij}) x_i x_j \\
 &= \frac{1}{2} \sum_i d_i x_i^2 + \frac{1}{2} \sum_j d_j x_j^2 - \sum_{ij} A_{ij} x_i x_j \\
 &= \frac{1}{2} \sum_{ij} [A_{ij} x_i^2 - 2A_{ij} x_i x_j + A_{ij} x_j^2]
 \end{aligned}$$

Here  $x_i$  and  $x_j$  represent diff. values of variable  $x$ .  
And not different dimensions.

arbitrary vector

across all 'i's and 'j's.

$$\cancel{x^T L x} = \frac{1}{2} \sum_{ij} A_{ij} (x_i - x_j)^2 \geq 0$$

always  $> 0$  1/0  
 $\therefore x^T L x \geq 0$   $x^T L x > 0$

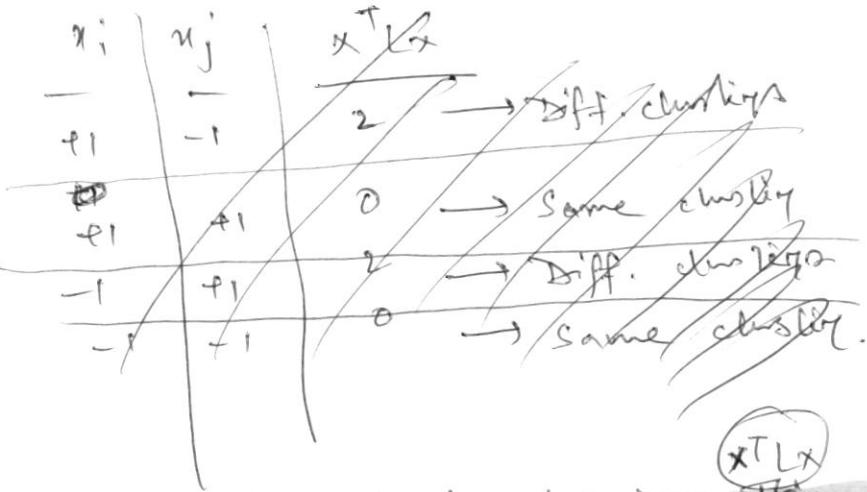
ApplicationApplication of  $x^T L x \geq 0$ ?CLUSTERING IN GRAPH  
NODES OF A GRAPH

Let's say, we want to cluster clustering a graph into two groups.  
Collective graph group based on similarity.

Note: - Graph clustering is not graph partitioning  
(dividing a graph into components based on connections)

(grouping nodes based on similarities and not based on edges)

How can  $x^T L x$  be useful in graph clustering?  
In  $x^T L x$ , let  $x$  be a vector that is (1) if part of cluster A and (2) if part of cluster B and 0 if  $x_i$  and  $x_j$  are in same clustering and 2 if they are in different clusters.



If we find the  $x$  that minimizes  $x^T L x$ , we will find clustering assignments that minimize cross-cluster edges

$n_i$	$n_j$	$A_{ij}$	$(x^T L x)_{ij}$
-	-	1	2
+1	-1	1	0
+1	+1	1	2
-1	+1	1	0
-1	-1	1	0
+1	-1	0	0
+1	+1	0	0
-1	+1	0	0
-1	-1	0	0

The same cluster if connected  
Two nodes in  
Or no connection b/w nodes

Max. value of  $(x^T L x)_{ij} = 2$

Max. Min. value of  $(x^T L x)_{ij} = 0$

$x^T L x$  will be minimized if 2 nodes that are connected lie in the same cluster or if there is no connection at all.

→ Instead of +1s and -1s, the values in  $x$  can be real, continuous values — embeddings of nodes. Then, minimizing  $x^T L x$  could be minimizing the distance b/w connected nodes. If there is no connection b/w  $n_i$  and  $n_j$ , we don't care about the distance as column

### How To Cluster Nodes

- Find  $x$  that minimizes  $x^T L x$ .  $x$  will contain the mapping of nodes to clusters.

If  $n_i$  and  $n_j$  are real, continuous values  $(n_i - n_j)^2$  represents the distance b/w them

the norm, comparison of by  
 along vector eigen vector. And, see count  
 smaller than eigen value, more the comparison  
 eigen-value with value at the most

To a graph (e.g., L), eigen vector with smallest  
 vector

- Rayleigh-Ritz Method

i.e. it always all ratio are 0  
 the solid all ratio are 0  
 This will avoid the continue case  
 should be normalized - able

example of this condition  
 in order to find sign  
 compare all values of in being 0

(load, no + & the point)

1) centre of mass about the origin:  $\sum x_i = 0$

This part (approximate)

This part

Addendum: constraints

constraint, we need to put additional  
 vector, to avoid getting  $x$ , that is an all-equal constraint  
 to avoid getting  $x$ , that is an all-equal constraint

So,  $Lx = 0 \quad (x = 0)$

(i.e., as eigen vector ( $c_1, c_2, \dots, c_n$ ) with value  $x^T Lx$  to 0)

A constraint vector ( $c_1, c_2, \dots, c_n$ ) with value  $x^T Lx$  to 0

with eigen value 0, surely, multiply  $x^T Lx$   
 A constraint vector ( $c_1, c_2, \dots, c_n$ ) is an eigen vector of

we need additional constraint between

15.14

v5.15

But, the smallest eigen vector corresponding to the smallest eigen value  $\lambda_0$  is an all-equi-constant vector

an all-equi-constant(1) vector  $\leftarrow$  We are trying to avoid this case (v5.14)

$$\lambda = 0 \Rightarrow x = 1$$

### Laplacian(L) Properties.

- 1) Real - Valued Elements
- 2) Symmetric:  $i$ th row =  $i$ th column  $\Rightarrow L = L^T$
- 3) Non-negative Eigen values ( $\lambda$ )
- 4) Orthogonal Eigen vectors ( $x$ )
- 5) Sum of elements of a row = 0  $\Rightarrow \lambda_1 = 0.1$   
 $\therefore 1$  is an Eigen vector of  $L$  with  $\lambda = 0$
- 6)  $x^T L x > 0$
- 7) Multiplying 'L' by The embeddings of a node is equiv. to

So, we use the ~~smallest~~ Eigen vector with the 2nd smallest Eigen-value. This soln. meets the 2 additional constraints ①  $x \perp 1$  and ② Normalizable to 1

② Normalizable to 1  
 i.e. not an all-zero constant vector

This Eigen vector is called "Fiedler Vector"

For example, the following Friedler vector can be converted to cluster assignment as follows:-

$$X = \begin{bmatrix} -0.5 \\ 0.2 \\ -0.1 \\ 0.3 \end{bmatrix} \rightarrow \begin{bmatrix} 1 \\ +1 \\ -1 \\ +1 \end{bmatrix}$$

Embeddings  $> 0$  are assigned +1 cluster

# How To Custer Nodes Of A Graph

▷ Construct Laplacian of a graph.

2) Get all its Eigen vectors and Eigen values

3) Get the Eigen vector(s) corresponding to 2nd smallest

4) Use this vector to draw.

4) Use the elements of this vector to classify the corresponding nodes of the graph.

→ So, the 2nd smallest eigen-value's Eigen Vectors  
is the  $X$  in  $X^T L X$  and its elements act as  
the embeddings of the corresponding nodes

$\{$   $i$ th element of the Fiedler vector in the embedding of the  $i$ th node of the graph

embedding of  
A  $k$ -dimensional embedding for a node can be constructed from the corresponding element of the 1st  $k$  (and smallest to the  $(k+1)^{\text{th}}$  smallest eigen-value) eigen vectors.

VS.17

- Just like as the 2nd smallest Eigen vector eigen-vector can be used to assign cluster to the nodes of ~~the~~ a graph, ~~to~~ ~~successively~~ successive larger eigen-valued Eigen vectors can also be used to assign clusters to the nodes of the graph.
- The cluster assigned by using the 2nd smallest  $\lambda$  is based on coarse properties.
- As we go on using larger  $\lambda$  ~~for~~,
- Clustering using larger  $\lambda$  is a
- As we ~~go on~~ go on to use larger and larger  $\lambda$ , the clustering is based on ~~more~~ finer and finer properties.
- As we saw that the elements of an Eigen vector can be ~~assigned~~ assigned as embeddings to the corresponding ~~elements~~ nodes of the graph.
  - Elements of larger eigen-valued Eigen vectors represent ~~properties~~ node features or embeddings that are more ~~fine~~ fine-grained.

In  $x^T L x$ , if  $x$  be a normalized eigenvector,

then,  $x^T L x = x^T \lambda x = \lambda x^T x = \lambda$  ( $x^T x$  is normalized)

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

$$= 1 \cdot 1 + 1 \cdot 1 + 1 \cdot 1$$

Because a Laplacian always has non-negative eigen values,

$$x^T L x = \lambda \geq 0 \quad \text{for normalized } x$$

~~If~~ Since eigen values of a Laplacian is always non-negative; they can be treated as frequencies because frequencies are always non-negative.

Small eigen value  $\Rightarrow$  low frequency.

Large eigen value  $\Rightarrow$  high frequency.

Elements of large Eigen values - Eigen vectors represent node-features that have higher frequency. Larger is the Eigen value, higher is the frequency of node features.

~~V5.17~~ ~~V5.19~~

## Properties of Laplacian ( $L$ )

1) Real-valued elements

2) Symmetric :  $i$ th row =  $i$ th column  $\Rightarrow L = L^T$

3) Non-negative Eigen values ( $\lambda$ )

4) Orthogonal Eigen vectors  $\Rightarrow \lambda_1 \cdot \lambda_2 = 0$

5) Sum of elements of a row = 0  $\Rightarrow \lambda_1 = 0.1$   
 $\Rightarrow 1$  is an Eigen-vector of  $L$

6) 1 is the Eigen vector with the smallest  $\lambda$  ( $= 0$ )

7)  $x^T L x \geq 0$  (semi-positive definite)

8) Multiplying ' $L$ ' with the features/embeddings of a node has the effect of getting the difference of a node's signal with its neighbours.

(14) Each Eigen vector  
of  $L$  has  $N$  elements  
 $N = \text{no. of nodes}$

9) Eigen vectors of  $L$  give us useful embeddings  
of the nodes

(13) No. of Eigen  
vectors of  $L$   
= No. of nodes  
in the graph

10) Smaller Eigen vectors with smaller  $\lambda$  give us coarse properties/features/embeddings of the nodes.

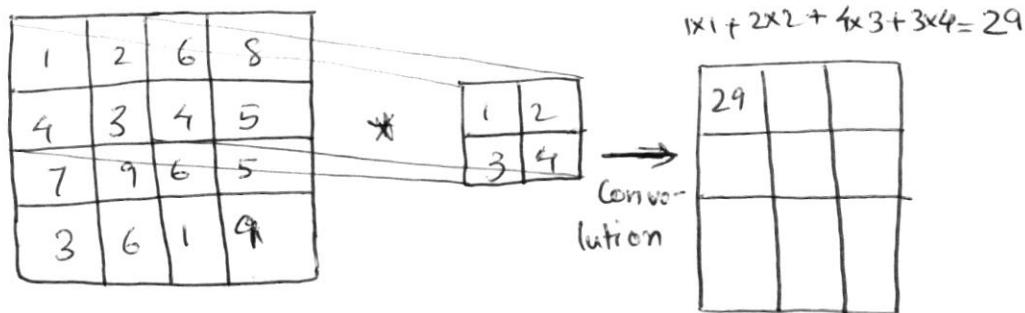
11) Eigen vectors with larger  $\lambda$  give us finer

12)  $L = U \Lambda U^T$ ,  $\Lambda = \text{diag}([\lambda_0, \dots, \lambda_{N-1}])$ ,  $U = \text{matrix of Eigen-vectors of } L \text{ stacked as columns}$

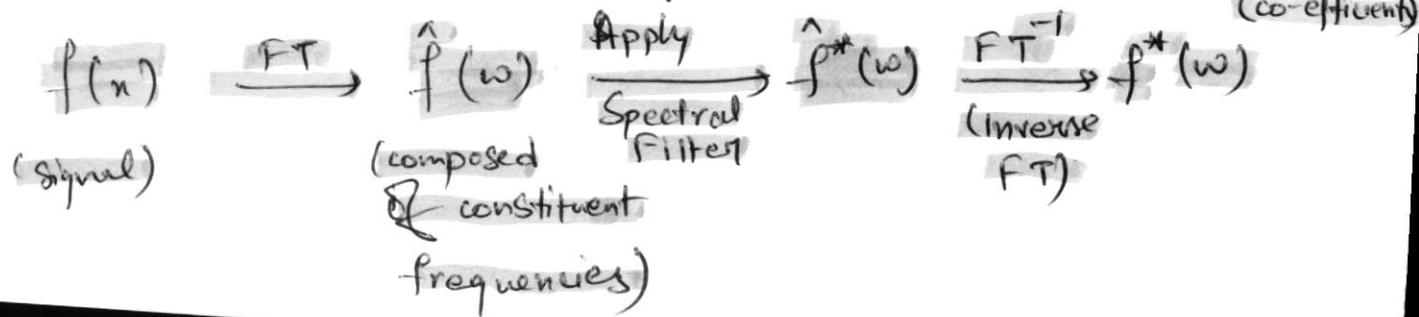
16.1

## Fourier Transforms, Wavelets and Spectral Convolutional Networks

- In Convolution Neural Networks (CNNs), the convolution operation is done by scanning the image-grid of pixels with a ~~small~~ filter - a small grid of values. The filter scans the entirety of the image.
- A significant amount of work was done ~~in~~ on how to perform convolution on graphs - like, how to even define a filter ~~+~~  
→ ~~convolve~~ Since the no. of neighbours is different for different nodes unlike in the case of image



- Spectral methods ~~do~~ deal with frequencies. Frequencies come from Fourier Transform (F.T). ~~Fourier~~ FT decomposes a ~~Signal~~ <sup>Signal</sup> into its component frequencies.
- Spectral filters are a ~~scalar~~ function ~~that~~ operate on signals and diminish/amplify ~~certain~~ different component frequencies of the signal ~~to diff by different factors~~ (co-efficients)



They have infinite elements  
 represent the frequency of said bins since they are continuous,  
 In continuous  $F_T(f(x))$  is the signal and  $e^{-j\omega x}$   
 multiplied and summed.  
 In Dot product, elements of the vector are  
 corresponding to the signal.  
 This tells how much of the frequency is present  
 frequency, the result will be a scalar (a coefficient)  
 a dot product (of inner product) in the signal and The  
 and the only vector to be the signal. If we take  
 frequency components of the vector to be a component frequency

Dot product of two vectors  $\Rightarrow$  how much one vector  
 along another vector

$F_T$  gives a set of coefficients Total. say how much of  
 each of these the frequency exist in a signal.

$$\begin{aligned}
 & \text{frequency} \\
 & \text{of the} \\
 & \text{continuous} \\
 & \text{fourier transform} \\
 & \text{is given by} \\
 & F_T \{ f(x) \} = \int_{-\infty}^{\infty} f(x) e^{-j\omega x} dx = \text{A scalar} \\
 & \text{in Space or time} \\
 & \text{dependent on} \\
 & \text{frequency} \\
 & \text{of the} \\
 & \text{continuous} \\
 & \text{fourier transform}
 \end{aligned}$$

1.6.2  
 1.7.1

6.3

Dot product is the sum of the products of the corresponding elements of 2 vectors  $\mathbf{x}_1$  and  $\mathbf{x}_2$

$$\boxed{\mathbf{x}_1 \cdot \mathbf{x}_2 = \sum_i x_{1i} x_{2i}}$$

Analogy :- Discrete domain Continuous Domain

$$\mathbf{x}_1 \rightarrow f(n)$$

$$\mathbf{x}_2 \rightarrow e^{-j\omega n}$$

$n \rightarrow \infty$  (infinite elements)

$$\sum_i x_{1i} x_{2i} \rightarrow \int_{-\infty}^{\infty} f(n) e^{-j\omega n} dn$$

$$\therefore FT \{f(n)\} = \hat{f}(\omega) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} f(n) e^{-j\omega n} dn$$

gives a co-efficient for a frequency  $\omega$ .

Continuous  
Inverse FT

$$FT^{-1} \{ \hat{f}(\omega) \} = f(n) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \hat{f}(\omega) e^{j\omega n} d\omega$$

adds up all the ~~for~~ component  
corresponding to a frequency  $\omega$  to obtain the ~~signal~~  
basis fun.  $\{e^{j\omega n}\}$  original signal.

$$\text{In } \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \hat{f}(\omega) e^{j\omega n} d\omega, \hat{f}(\omega) \text{ is projected}$$

onto a new set of  
functions that ~~are~~ are parame-  
terized by  $\omega$ :  $e^{j\omega n}$  ( $= \cos \omega n + j \sin \omega n$ )

16.4

$e^{i\omega n}$  is the eigen fn. of  $\frac{\partial^2}{\partial x^2}$

$$\therefore \frac{\partial^2}{\partial x^2} (e^{i\omega n}) = -\omega^2 e^{i\omega n}$$

(iMy,  $Lx = \lambda x$ )

Here,  $\frac{\partial^2}{\partial x^2}$  is 1-D Laplacian operator of  $e^{i\omega n}$   
because when you plug in the fn. in the  
operator, you get the same fn.  
scaled by a constant.

- ① The basis fns. ~~are~~ of the FT are eigen functions of 1-D Laplacian operator
- ② FT  $\rightarrow$  Expansion Expressing a signal in terms of its component frequencies.
- ③ Continuous FT  $\rightarrow$  Expressing a function in terms of eigen functions of the Laplacian operator
- ④ Graph FT  $\rightarrow$  Expressing in terms of eigen vectors of the Graph Laplacian, L

Continuous  $\rightarrow$  Discrete

Operator:  $\rightarrow \frac{\partial^2}{\partial x^2} \rightarrow L$

$e^{i\omega n}$   $\rightarrow x$

$$\boxed{\begin{array}{ccc} \frac{\partial^2}{\partial x^2} (e^{i\omega n}) & \rightarrow -\omega^2 (e^{i\omega n}) \\ Lx & \rightarrow \lambda x \end{array}}$$

W6.5

## Graph Fourier Transform

 $\lambda_0 = 0$ 

**GRAPH - LEVEL**

$$L X_e = \lambda_e X_e \quad (\lambda : 1, 2, 3, \dots)$$

$\lambda_0 < \lambda_1 < \lambda_2 < \lambda_3 < \dots$

**Graph FT:**

$$f(\ell) = \langle X_\ell, f \rangle = \sum_{n=1}^N X_\ell^*(n) f(n)$$

$\langle , \rangle$  inner product

$f(\ell)$  is the scalar co-efficient for the eigen vector  $X_\ell$ .

A scalar representing representation at the graph level

Summed across all the nodes of the graph

eigen vector of  $L$  of the graph

features of a node

Inverse Graph FT:

$$f(n) = \langle \hat{f}, X_\ell \rangle = \sum_{\ell=0}^{N-1} \hat{f}(\ell) X_\ell(n)$$

**NODE - LEVEL**

Similar to a signal, a graph can be expressed as the sum of the product of eigen vectors and corresponding co-efficients ( $f(\ell)$ )

Summed across all eigen vectors and  $(X_\ell)$   
 $F^T_S (\hat{f}(n))$  of the graph

A scalar Graph-level

$f(\ell)$  tells us how much of the frequency/Eigen vector ' $X_\ell$ ' is present in the entire graph

$X_\ell(n)$ : -  $n$ th element of the Eigen-vector

No. of elements in it = no. of nodes

$X_\ell(n)$ : It corresponds to the  $n$ th node

$f(n)$ : -  $n$ -th element of the graph.  
 $\hat{f}(n)$ : Feature vector of  $n$ th node  
 This is our signal

$f(n) = \text{a scalar}$   
 $\hat{f}(n) = \text{a scalar}$

W6.5

\* \* \* But, a node has a vector of features.

v6.6

$X_l$  = Eigen vector: constituent frequency

n<sup>th</sup> element of  $X_l$  corresponds to feature of n<sup>th</sup> node ~~of~~

Feature Vector of  
n<sup>th</sup> node

$$\begin{bmatrix} x \end{bmatrix}$$

n<sup>th</sup> element  
of each ~~eigen~~ Eigen vector

$$\begin{bmatrix} X_2(n) \\ X_3(n) \\ X_4(n) \\ \vdots \\ X_m(n) \end{bmatrix}$$

$\hat{f}(l)$  tells us how much of the n<sup>th</sup> element of the l<sup>th</sup> Eigen Vector  $X_l$  is present in the ~~n<sup>th</sup>~~ feature of the n<sup>th</sup> ~~node~~ node

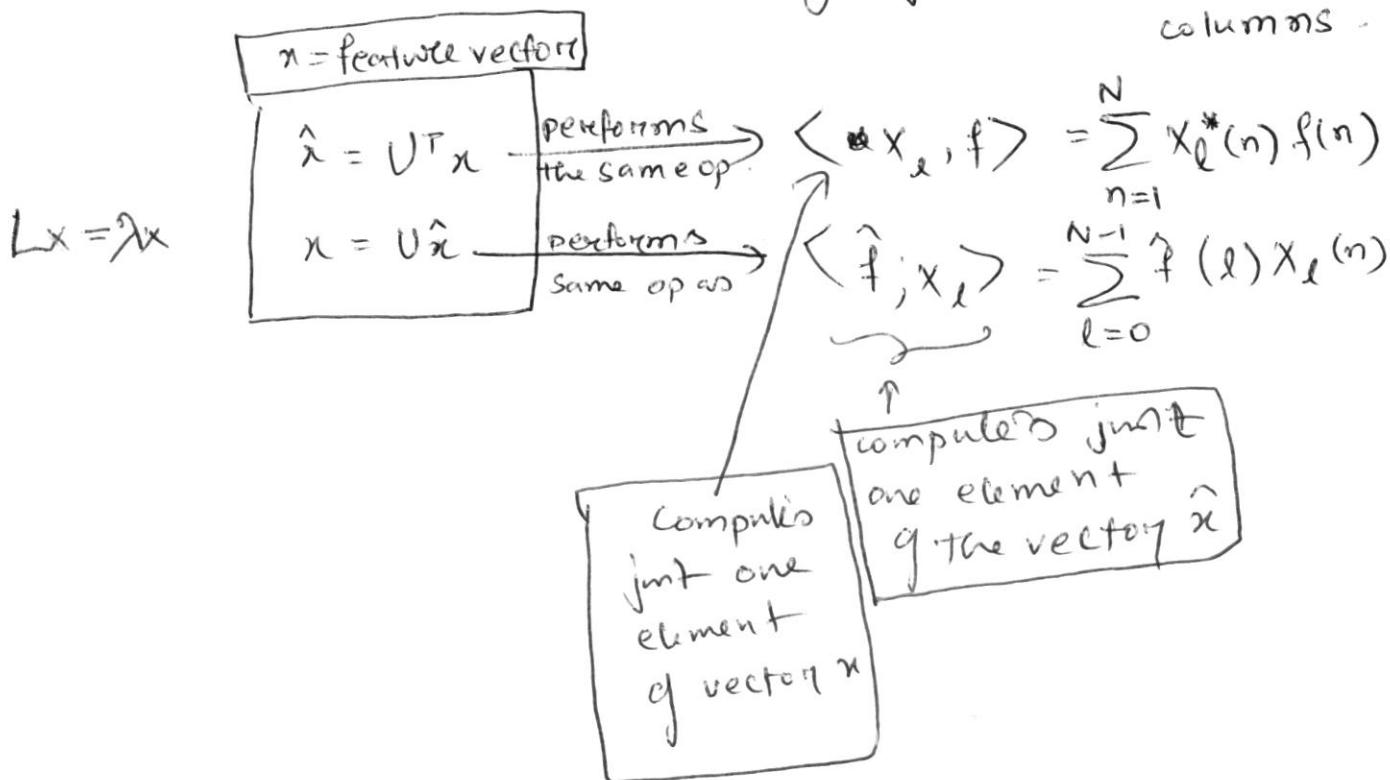
To get the total feature of the n<sup>th</sup> node, sum up the ~~n<sup>th</sup>~~ products of the n<sup>th</sup> elements of the all the Eigen vectors and their corresponding Fourier coefficients. That is an inner product. ~~of~~

v6.7

# Graph Fourier Transform in Matrix Form.

$$L = U \Lambda U^T, \quad \Lambda = \text{diag}([\lambda_0, \dots, \lambda_{N-1}])$$

$U$  = Matrix of Eigen vectors of  $L$ , stacked as columns.



In (v 6.1),  $f(n) \xrightarrow{\text{FT}} \hat{f}(\omega) \xrightarrow{\substack{\text{Apply} \\ \text{spectral} \\ \text{Filter}}} \hat{f}^*(\omega) \xrightarrow{\text{FT}^{-1}} f^*(n)$

In case of graph,

$$f(n) \xrightarrow[\text{(FT)}]{U^T x} \hat{f}(\omega) \xrightarrow[\text{& filter}]{\substack{\text{Apply} \\ g_\theta^*(\omega)}} \hat{f}^*(\omega) \xrightarrow[\text{(FT}^{-1}\text{)}]{U \hat{x}} f^*(n)$$

$g_\theta^*(\omega)$  is the filter spectral filter

We apply this filter for performing spectral convolution

Step 1:

$$x^* = U g_\theta^*(\omega) U^T x$$

This is how convolution is performed on graphs

V 6.8

$$\mathbf{x}^* = \mathbf{U} \mathbf{f}_0(\Lambda) \mathbf{U}^T \mathbf{x}$$

↓  
 FT  
 ↓  
 Apply filter  
 ↓  
 FT<sup>-1</sup>

It's expensive to compute I use all the eigen vectors of L

[FT is mapping a fn. onto a set of basis functions]

Issues with Graph Spectral Filtering Problems:

→ If N = no. of nodes in the graph, no. of Eigen vectors = N, length of each Eigen vector = N. So, U is NxN big. For large graphs, U will be very big and hence, computationally expensive to use.

→ In calculating  $\mathbf{x}^*$  by using the above formulae, the FT co-efficient calculation ( $= \mathbf{U}^T \mathbf{x} \circledast, = \langle \mathbf{x}, f \rangle = \sum_{n=1}^N x_n^* f(n)$ )

involves all the nodes of the graph. When these co-efficients are used to compute  $\mathbf{FT}^{-1}$ , all the other nodes of the graph influence this computation.

The re-construction of the feature vector of a node. So, the process of taking FT, applying spectral filter and then taking FT<sup>-1</sup> is highly non-local. We prefer only the immediate neighbourhood nodes to influence the node's feature computation.

These 2 problems get filters solved by using graph wavelets in FT computation.

Graph wavelets can be approximated by low-order polynomials of the Laplacian and as an added benefit, it avoids the expensive eigen decomposition computation.

6.9

The problems filters  
Graph wavelets can be approximated by low-order  
(Chebyshev polynomials of the Laplacian to:-)

① Localize filter - avoiding the influence of distant nodes (Problem 1)

② expensive  
avoid eigen decomposition computation (Problem 2)

(Eigen decomposition of  $L = U \Lambda U^T$ )

$$y = g_\theta(L)x = g_\theta(U\Lambda U^T)x = U g_\theta(\Lambda) U^T x$$

Instead of doing this,  $g_\theta(L) = \sum_k \theta_k T_k(L)$

~~$x^* = g_\theta(L)x$~~

~~$x^* = g_\theta(L)x \approx \sum_{k=0}^K \theta_k T_k(L)x$~~

~~$\Rightarrow x^* \approx (\theta_0 + \theta_1 \tilde{L})x$  for  $K=1$~~

~~Filter~~

~~① Localizes the operation to 1-hop neighbourhoods~~  
~~② Simplifies computation by avoiding  $U$  and  $U^T$~~

$\theta_0$  and  $\theta_1$  are filter parameters that need to be learnt. They are learnt via BackPropagation

Final  
Graph  
Convolutional  
Operation

$$x^* \approx (\theta_0 I + \theta_1 \tilde{L})x$$

Computationally Simple

$I$  = Identity matrix

$\theta_0 K=1$   
 $\theta_1$  It's 1st order approximation

So, GNN is ~~is~~ So, the source of GNN lies in :-

A 1st order approximation of wavelet filters on graphs ( $K=1$ )

### Some more approximations:-

①  $\tilde{Z}$  represents scaled, normalized  $Z$

$$L = I - D^{-1/2} A D^{-1/2} \quad (\text{D-Degree matrix})$$

(Normalized) A - Adjacency matrix

$$\tilde{L} = \frac{2}{\lambda_{\max}} - I \quad (\lambda_{\max} \text{ is largest } \lambda)$$

$$= I - D^{-\frac{1}{2}} A D^{-\frac{1}{2}} - I \} f_{\alpha \gamma} \}_{\max}$$

$$\Rightarrow T = - D^{-\frac{1}{2}} A D^{-\frac{1}{2}}$$

$$\begin{aligned} S_0, x^* = g_0 x &\approx (\theta_0 I + \theta_1 \tilde{L})x \\ &\approx (\theta_0 I - \theta_1 D^{-1/2} A D^{-1/2})x \end{aligned}$$

$$\textcircled{2} \quad \theta_0 = -\theta_1 = \theta$$

$\theta = -\theta_1 = \theta$   
~~This~~ This allows reduction of in no. of parameters  
~~in the final expression~~

Also, it ~~creates~~ simplifies / beautifies the expression  
of  $x^6$  into 3 distinct factors:-

$$x^* = g_\theta x \approx \theta(I + D^{-1/2} A D^{-1/2})x$$

↑ parameter      ↑ info regarding graph      ↑ node feature vector

③ The factor  $(I + D^{-1/2} A D^{-1/2})$  which is a matrix has Eigen values ~~not~~ in  $[0, 2]$ . So, repeated applications of this matrix will blow up the scale. [WHY WE NEED REPEATED APPLICATIONS]

So, the following "renormalization" trick is performed.

$$\tilde{A} = A + I \quad (\equiv \text{Adding a self-loop edge to each node})$$

$$\bar{d}_{ii} = \sum_j \bar{A}_{ij} = D_{ii} + 1$$

$$\left( \mathcal{F} + \bar{D}^{\frac{1}{2}} \mathcal{A}_{AD}^{-\frac{1}{2}} \right) \rightarrow \bar{D}^{-\frac{1}{2}} \bar{\mathcal{A}} \bar{D}^{-\frac{1}{2}}$$

$$\text{So, } n^* = g_{\theta}^* x \approx \underbrace{\theta}_{\substack{\uparrow \\ \text{Filter}}} (\underbrace{\tilde{D}^{-1/2} \tilde{A} \tilde{D}^{-1/2}}_{\substack{\uparrow \\ \text{Graph}}}) \underbrace{n}_{\substack{\uparrow \\ \text{Feature}}}$$

✓72

## GCN Final Form

$$\begin{aligned} Z &= \underbrace{D^{-\frac{1}{2}} A D^{-\frac{1}{2}}}_{\text{Graph}} X \Theta \\ &= \underbrace{\hat{A} X \Theta}_{\text{Message Passing Operation}} \\ &\quad \uparrow \text{Graph} \quad \uparrow \text{Filter} \\ &\quad \underbrace{\text{Features}}_{\text{Features}} \end{aligned}$$

$\hat{A}X$  :-

[When we multiply a set of features by the adjacency matrix  $A$ , what we ~~are doing~~ do essentially is summing up ~~all~~ the features ~~of~~ of all its connections (nodes).]

[In  $Z = \hat{A}X\Theta$ , we are multiplying the features with the normalized adjacency matrix  $\hat{A}$  (which includes self-edges to every node). So, ~~it's~~  $\hat{A}X$  is equivalent to summing up the features of its neighbours as well as its own features (because of the self-edge).]

## Multi-Layered GCN

