

# Graph Convolutional Networks

## lecture-7

Course on Graph Neural Networks (Winter Term 21/22)

Prof. Dr. Holger Giese ([holger.giese@hpi.uni-potsdam.de](mailto:holger.giese@hpi.uni-potsdam.de))

Christian Medeiros Adriano ([christian.adriano@hpi.de](mailto:christian.adriano@hpi.de)) - “Chris”

Matthias Barkowski ([matthias.barkowski@hpi.de](mailto:matthias.barkowski@hpi.de))

- Metrics
- Classification
- Random Walk
- Node and Graphs Embeddings
- Graph Structure Learning

## **Why do we need Graph Neural Networks?**

We need a way to scale without having to make large random walks in the graph

Example from Industry

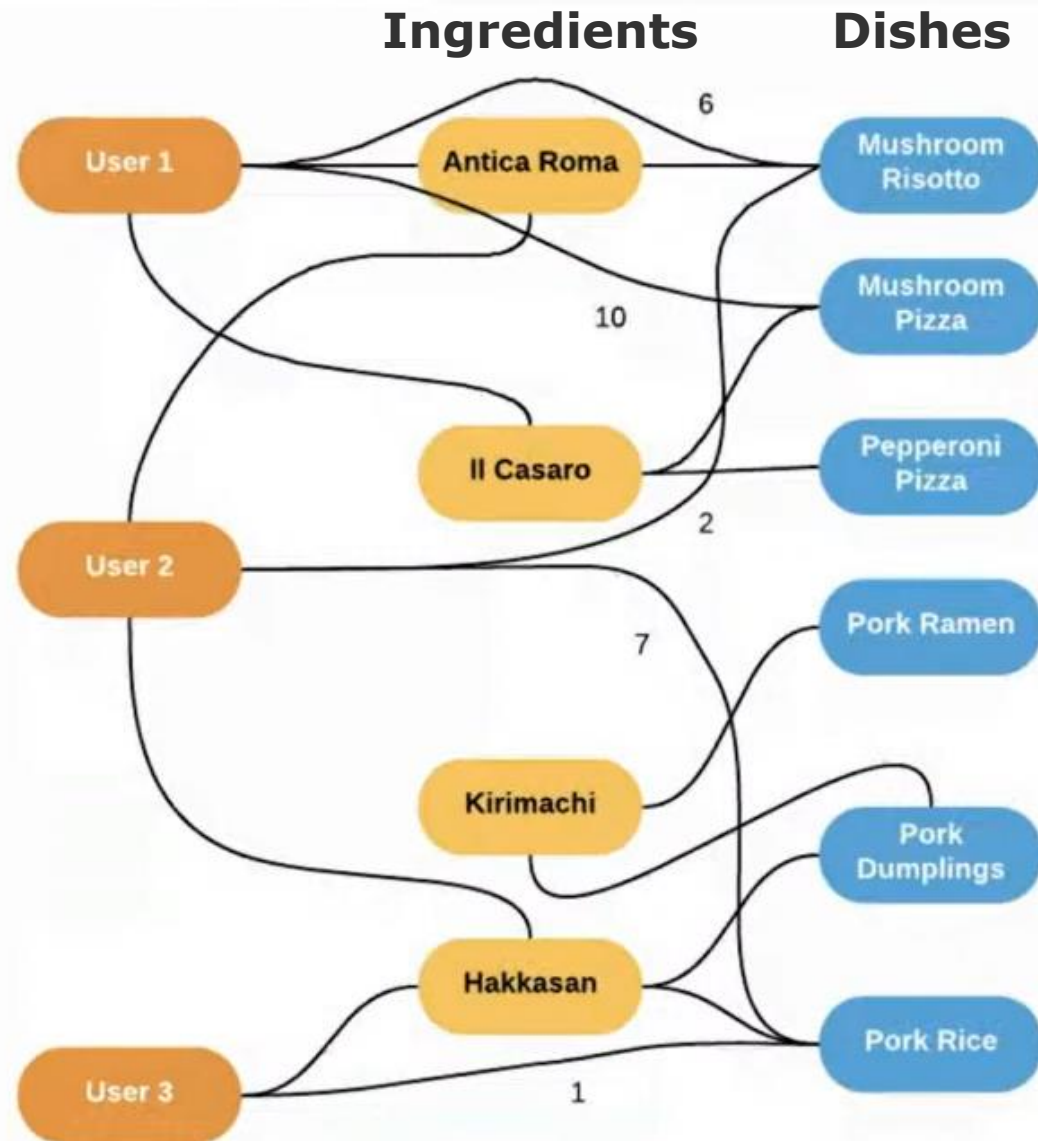
Graph Convolutional Networks (GCN)

Relational GCN

Mathematical background

- Laplacian
- Chebyshev Approximation
- Softmax

# Motivation – Uber Eats



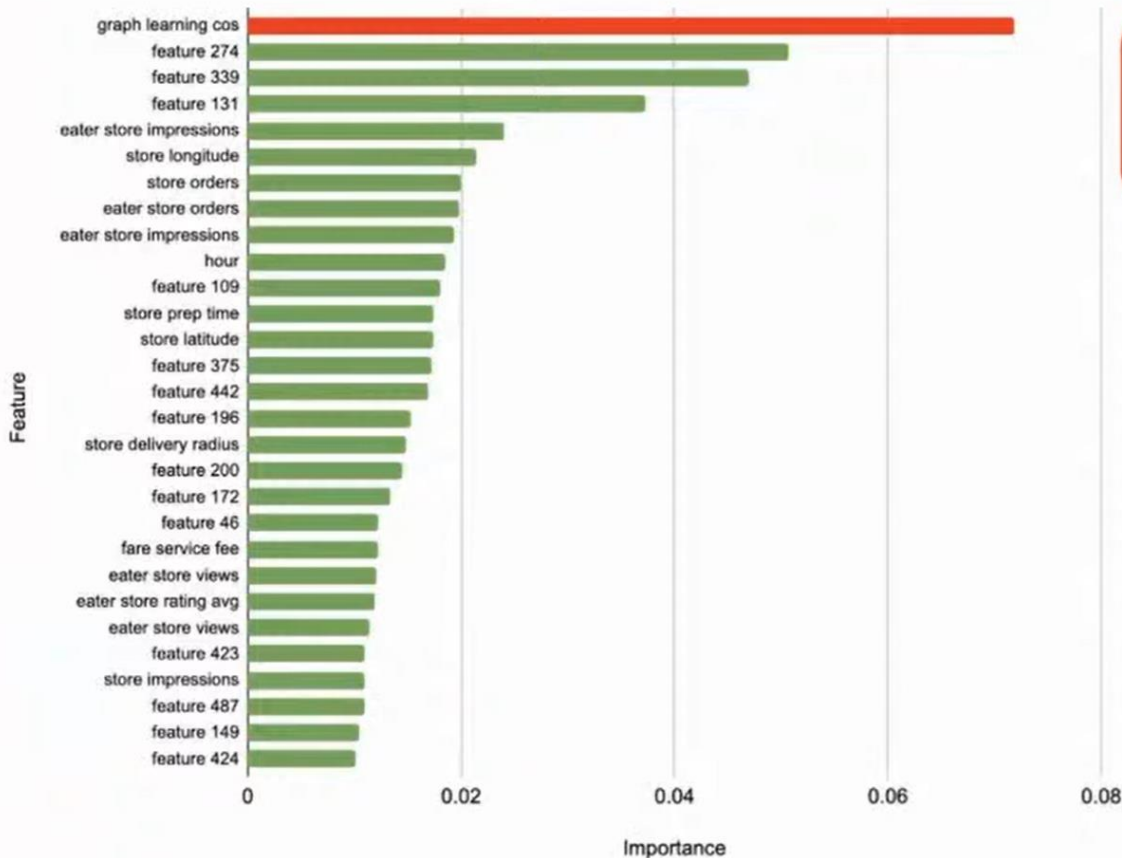
## Offline evaluation

Trained the downstream Personalized Ranking Model using graph node embeddings

~12% improvement in test AUC over previous production model

Model	Test AUC
Previous production model	0.784
With graph embeddings	<b>0.877</b>

## Feature Importance



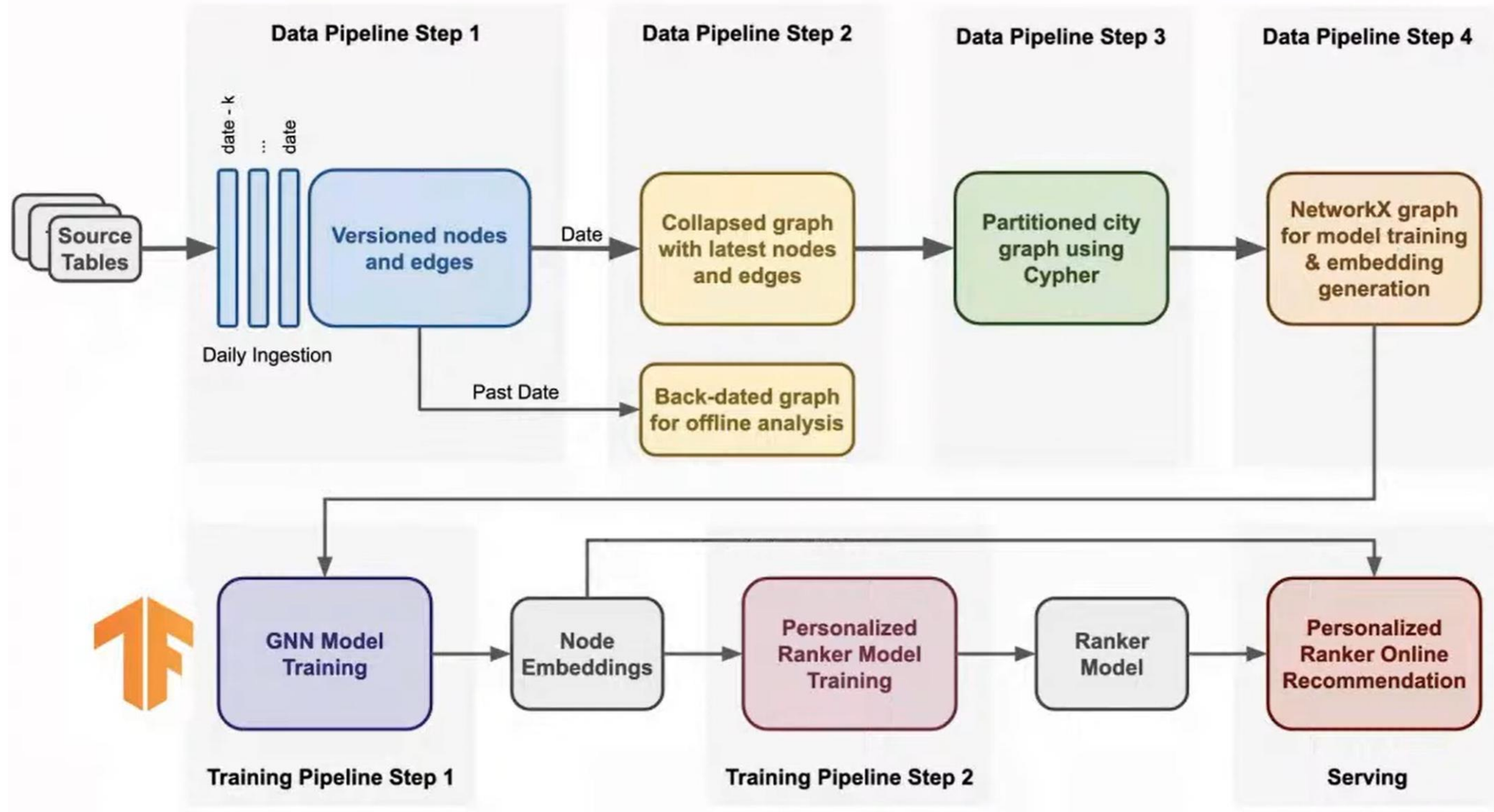
Graph learning cosine similarity is the top feature in the model

## Online evaluation

Ran a A/B test of the Recommended Dishes Carousel in San Francisco

**Significant** uplift in Click-Through Rate with respect to the previous production model

**Conclusion:** Dish Recommendations with graph learning features are live in San Francisco, soon everywhere else



# Bipartite Graph for Dish Recommendation

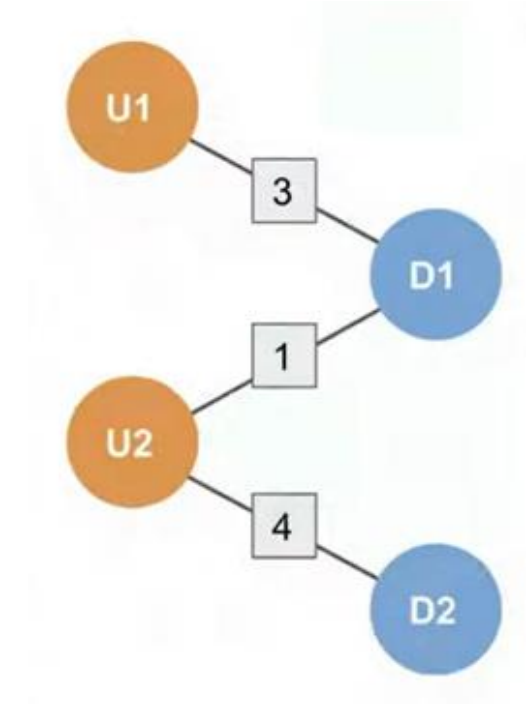
$U_i$  users

$D_i$  Dishes

Users connect to dishes by number of times they order a dish  $D_i$

New users and new dishes are added every day

Each node has a feature (e.g., word2vec)



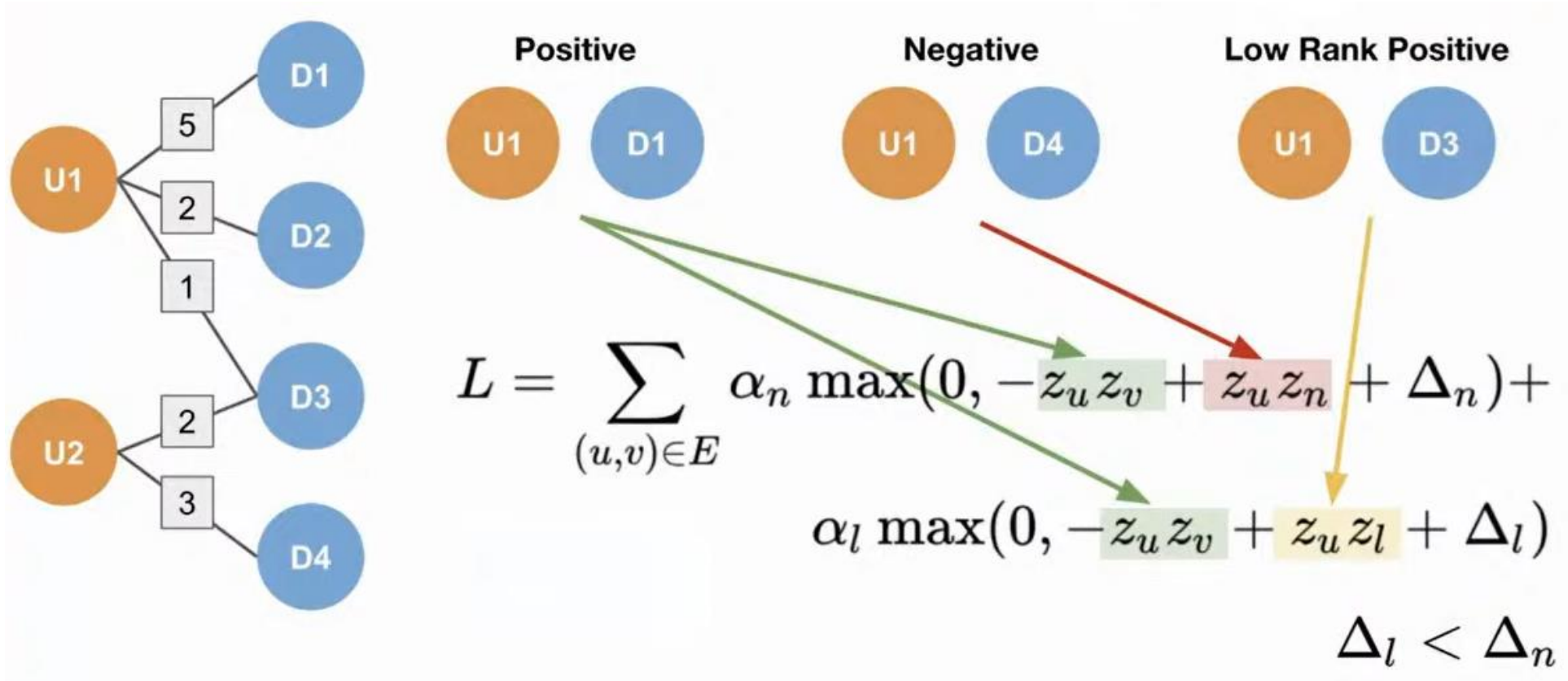
For dish recommendation, the goal is to rank dishes for each user.

$$L = \sum_{(u,v) \in E} \max(0, -\underbrace{z_u z_v}_{\substack{\uparrow \\ \text{positive} \\ \text{pair}}} + \underbrace{z_u z_n}_{\substack{\uparrow \\ \text{negative} \\ \text{sample}}} + \underbrace{\Delta}_{\substack{\uparrow \\ \text{margin}}})$$

Goal is for the positive pair to surpass the negative pair

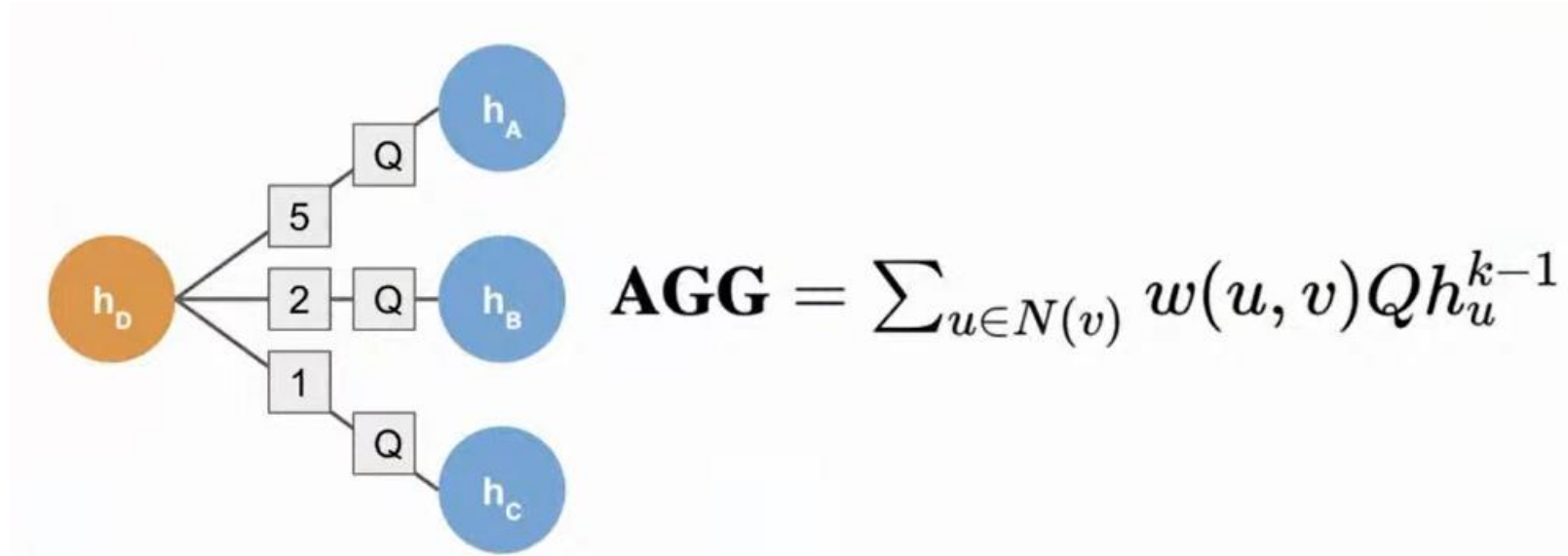


# Loss with Low Rank Positives



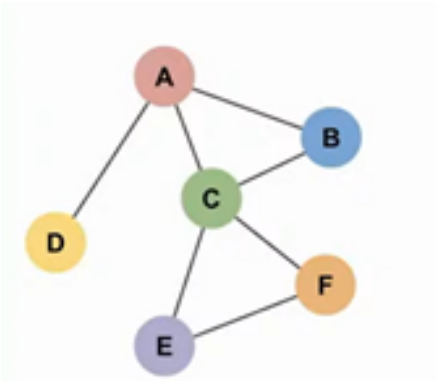
# Weighted pool aggregation

Aggregate neighborhood based on the edge weights

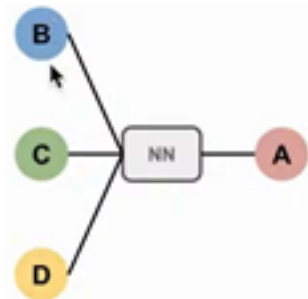


Obtain a node representation by using a neural network to aggregate the information from the neighbors.

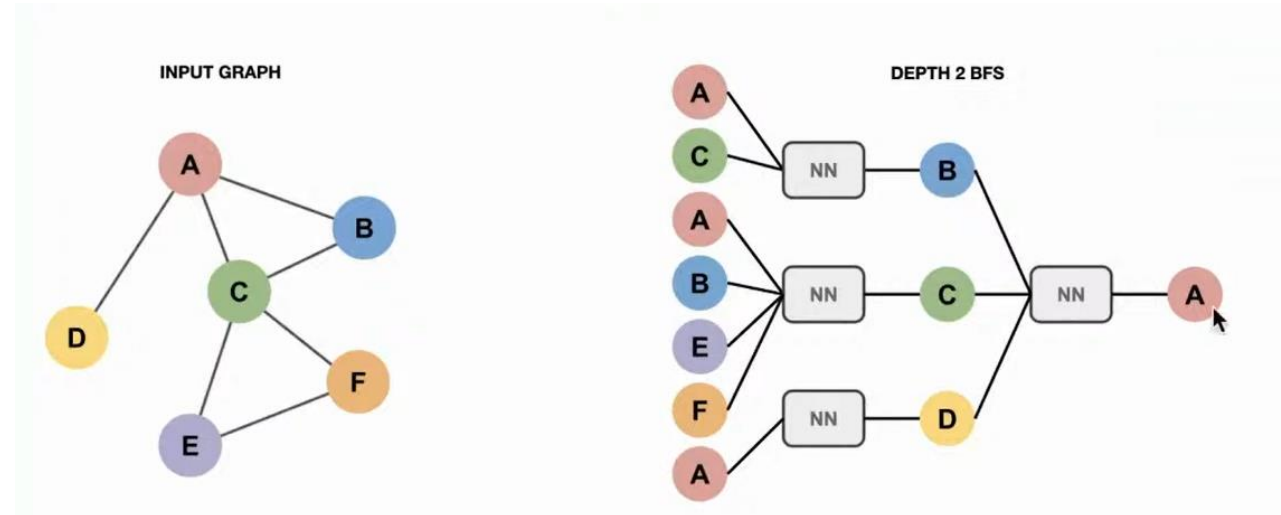
Input graph



Output graph  
(depth =2 BFS)



**Do it recursively**



# Solution for Recursion = Message Passing

Every node has its own embedding at every layer

Each layer is how deep we go in the network

Input feature  $x_u$  has layer K embeddings  $h_u^k$

Layer-0 has the features of each node

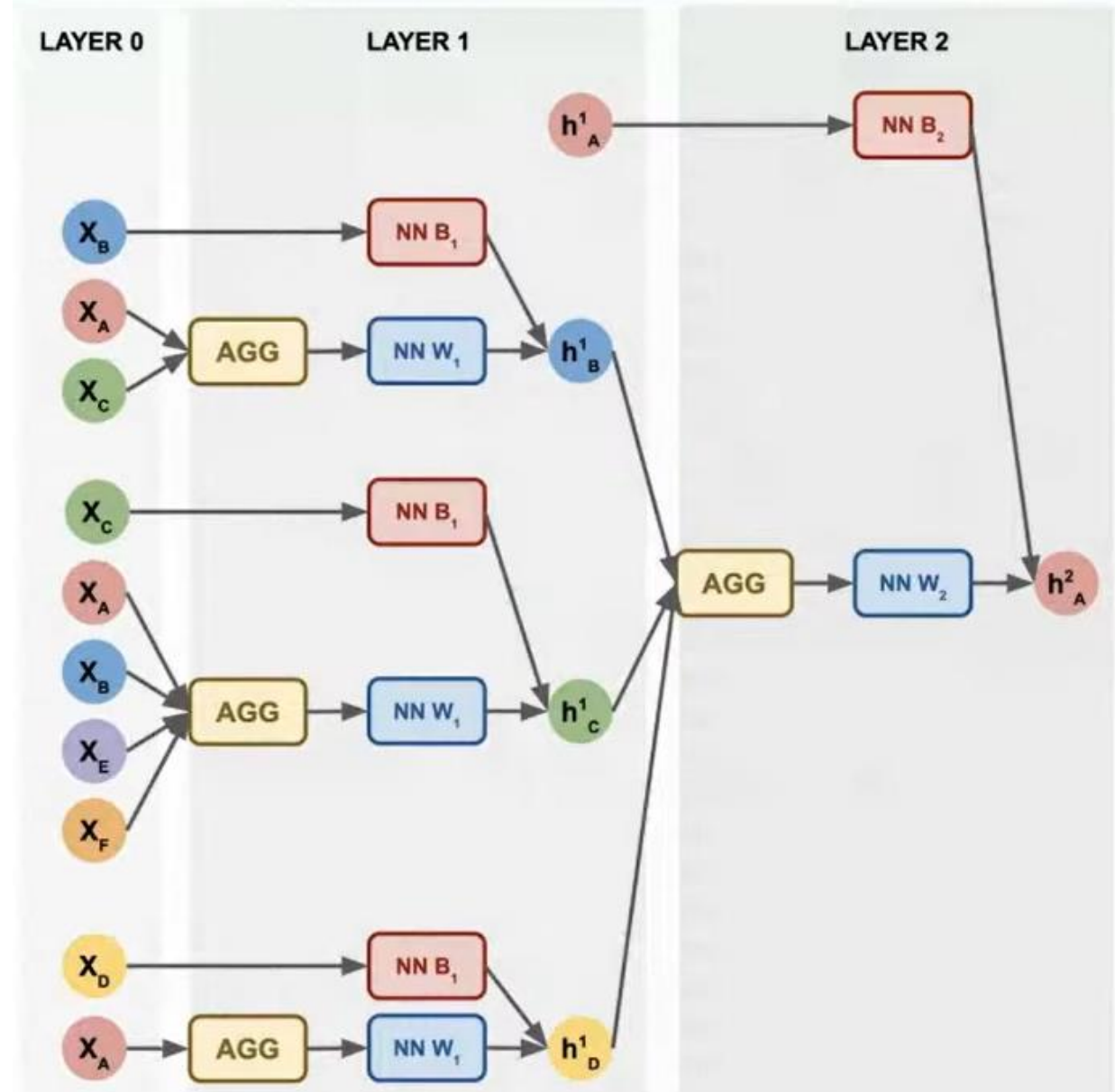
$AGG$  = aggregation function

$NN B_k$  = network self-embedding embedding

$NN W_1$  = network 1

$NN W_2$  = network 2

$h_u^k$  = computed representations



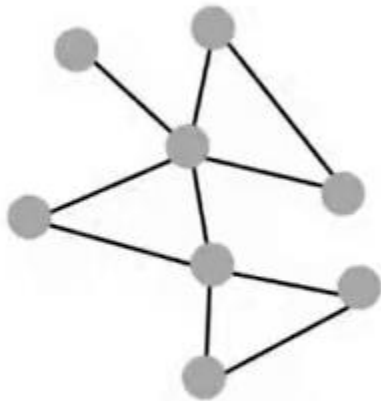
Message Passing in GNN allows to scale because:

1. We do not depend on the number of nodes in graph, but the number of neurons that we decide to use to compute de embeddings at each layer.
2. We can limit the number of neurons because all these nodes share the same network at each layer
3. This is the scalability that we do not have in the standard methods like Node2Vec, for instance.

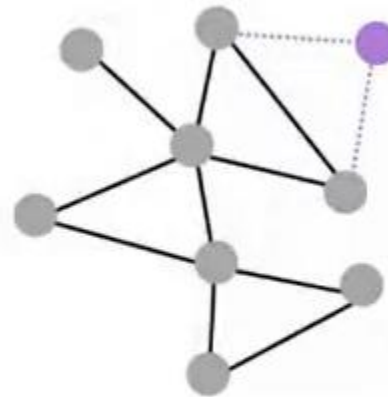
In many applications new nodes are constantly added to the graph

This would require constant retraining to update the embeddings

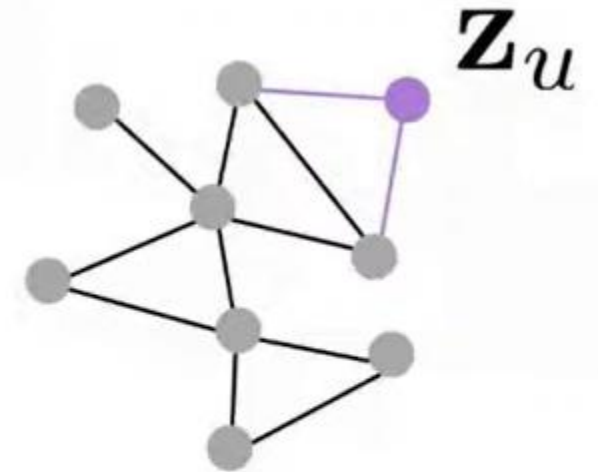
While this is difficult to do with traditional methods, GNN allows to retrain on the part affected



**train with snapshot**



**new node arrives**



**generate embedding  
for new node**

- Each node can have one or multiple dimensions (depends on the features that we choose)
- In a 2-layer GCN, each node pull the information from their neighbors two times.
- The number of GCN layer tells how far the signal could travel.
  - In a 2-layer GCN, the signal travels to a maximum to a two hops from the node.

# Steps

1. Message propagation
2. Feature aggregation
3. Encoding

## Label Propagation

✉ = Label

✉ = (fraud)

*label smoothing*

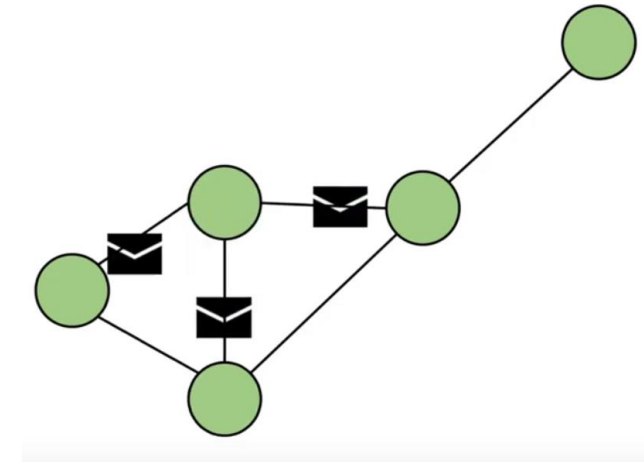
## GCN

✉ = Input Features

✉ = 

US	VN	IR
----	----	----

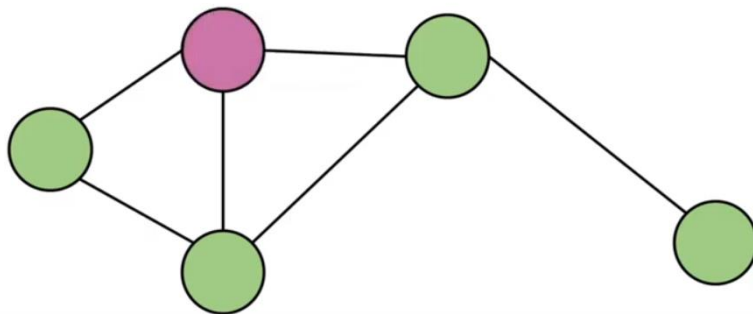
*feature smoothing*



Average( 

2.0	1.0	1.0
1.0	0.0	3.0

 )

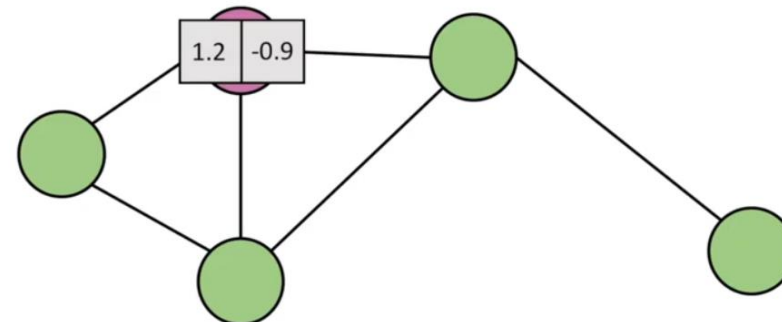


Average( 

2.0	1.0	1.0
1.0	0.0	3.0

 ) = 

0.7	0.3	0.7
-----	-----	-----





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

, where:

$l$  is a layer

$\sigma$  is the sigmoid function (adds the non-linearity)

$N_i$  are the neighbors of node  $i$

$c_{ij}$  is normalizing constant

$h^l$  is the value of the node in the previous layer

$W^l$  is a linear projection matrix for layer  $l$

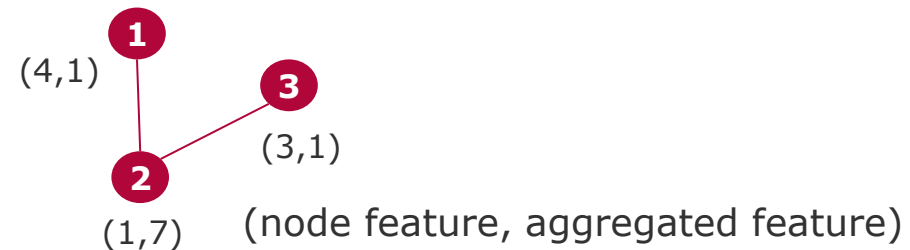
# Quick Math background based on Message Passing

$$\text{aggregate}(A_i, X)_i = A_i X = \sum_{j \in N} A_{i,j} X_j$$

This allow to compute the aggregate feature representation of the  $i_{th}$  node as vector-matrix product. Computationally, this weighted sum is obtained by using the Adjacency matrix to adjust or weight the features of a node by its he connectivity (edges or neighbors).

The contribution of the  $j_{th}$  node in the feature aggregation depends on value of the  $j_{th}$  column of the  $i_{th}$  row of  $A$ . Because  $A$  is the adjacency matrix, if this value is 1, then the  $j_{th}$  node is a neighbor of the  $i_{th}$  node.

$$\begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 3 \\ 1 \\ 4 \end{bmatrix} = \begin{bmatrix} (0)(3) + (1)(1) + (0)(4) \\ (1)(3) + (0)(1) + (1)(4) \\ (0)(3) + (1)(1) + (0)(4) \end{bmatrix} = \begin{bmatrix} 1 \\ 7 \\ 1 \end{bmatrix}$$



Hence, the contribution of each neighbor depends only on the connectivity specified in the adjacency matrix **A**

$$\text{aggregate}(A, X)_i = D^{-1} A_i X = \sum_{k \in N} D_{i,k}^{-1} = \sum_{j \in N} A_{i,j} X_j = \sum_{j \in N} \frac{A_{i,j}}{D_{i,i}} X_j$$

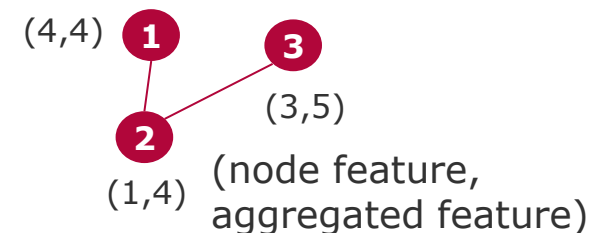
Where  $D$  is a diagonal matrix called the degree matrix, which has in its diagonal the degree of each node and zero in all other positions outside the diagonal. The inverse of  $D$  is still a diagonal matrix. This allow to remove first summation over  $D_{i,k}$  and use  $D_{i,i}$ .

When we divide the Adjacency matrix  $A$  by the degree matrix  $D$ , we obtain an average of the connections for each node, for example, if a node  $i$  has two connections, then it will get 0.5 in for each of the edges (or position in the  $A$  matrix). Hence, now we are multiplying the features by this averaged connection. **The consequence** is that neighbor nodes contribute individually (larger weight) to the aggregated feature of less connected nodes than of higher connect ones.

If one wants to include the features of the node itself, one can simply add 1 to the diagonal of the adjacency matrix  $A$ . Computationally this can be achieved by adding the identity matrix to the adjacency matrix, i.e.,  $\tilde{A} = A + I$

$$\begin{matrix} D^{-1} & A & I & X \end{matrix}$$

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1/2 & 0 \\ 0 & 0 & 1 \end{bmatrix} \left\{ \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} + \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \right\} \begin{bmatrix} 3 \\ 1 \\ 4 \end{bmatrix} = \begin{bmatrix} (1)(3) + (1)(1) + (0)(4) \\ (1/2)(3) + (1/2)(1) + (1/2)(4) \\ (0)(3) + (1)(1) + (1)(4) \end{bmatrix} = \begin{bmatrix} 4 \\ 4 \\ 5 \end{bmatrix}$$



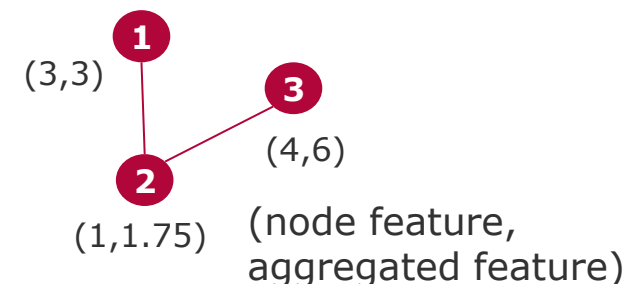
$$\text{aggregate}(A, X)_i = D^{-0.5} \tilde{A} D^{-0.5} X = \sum_{j \in N} \frac{1}{D_{i,i}^{0.5}} \tilde{A}_{i,j} \frac{1}{D_{j,j}^{0.5}} X_j =$$

This normalizes aggregated features so they remain roughly on the same scale as the input features. Another way to write is like the following:

$$\text{aggregate}(A, X)_i = \frac{1}{\sqrt{d_i d_j}} \tilde{A}_{i,j}$$

However, the spectral operation weighs neighbor in the weighted sum higher if they have a low-degree and vice versa. This is useful when low-degree neighbors carry more information than their high-degree counterparts.

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1/2 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 1 \\ 0 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1/2 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 3 \\ 1 \\ 4 \end{bmatrix} = \begin{bmatrix} (1)(3) + (0)(1) + (0)(4) \\ (1/2)(3) + (1/4)(1) + (0)(4) \\ (1/2)(3) + (1/2)(1) + (1)(4) \end{bmatrix} = \begin{bmatrix} 3 \\ 1.75 \\ 6 \end{bmatrix}$$



Hence, the contribution of each neighbor depends only on the connectivity specified in the adjacency matrix **A**

## Summing the weights of my neighbors

$\vec{X}$  is the feature vector, where each element  $i$  is a feature of a node  $i$

$A$  is the adjacency matrix

$$H = A * \vec{X}$$

## Average over the weights of my neighbors

$D$  is the node degree matrix, where  $D_{ii} = \sum_j A_{ij}$ ,

$$\tilde{A} = D^{-1} * A, D \text{ is the degree matrix}$$

## Scaling by the target node

$\tilde{A} = A * I$ ,  $I$  is the identity matrix,  $\tilde{A}$  adds the self-loops to the  $A$

$$\hat{A} = \tilde{D}^{-1/2} * \tilde{A} * \tilde{D}^{-1/2}$$

$$\hat{A}_{i,j} = \frac{1}{\sqrt{d_i d_j}} * \tilde{A}_{ij} \text{ scaling by the neighborhood sizes of source and root nodes}$$

This allows to have an adjacency matrix that if multiplied various times the computations will be stable, i.e., it will converge to a stationary value.

[https://github.com/christianadriano/blog\\_code/blob/master/gcn\\_numpy/message\\_passing.ipynb](https://github.com/christianadriano/blog_code/blob/master/gcn_numpy/message_passing.ipynb)



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

, where:

$W^l$  is a layer-specific trainable weight matrix.

$\sigma$  is an activation function, such as the  $\text{ReLU}() = \max(0; )$  or Sigmoid.

$H^l \in \mathbb{R}^{N \times D}$  is the matrix of activations in layer  $l$

$H^0 = \vec{X}$ , where  $X$  is the feature vector

Spectral Graph Convolution as the multiplication of a signal  $x$  by a filter  $g_\theta$

$$g_\theta x = U g_\theta U^T x, \text{ where:}$$

$U$  is the matrix of eigenvectors of the normalized graph Laplacian  $L$

$$L = I - \tilde{D}^{-\frac{1}{2}} * A * \tilde{D}^{-\frac{1}{2}} = U \Lambda U^T$$

$g_\theta$  is a function of the eigenvalues of the Laplacian, so  $g_\theta(\Lambda)$

However, computing this does not scale:

- Multiplication of  $U$  is  $O(N^2)$  and
- Eigen-decomposition of  $L$  is too expensive for large graphs

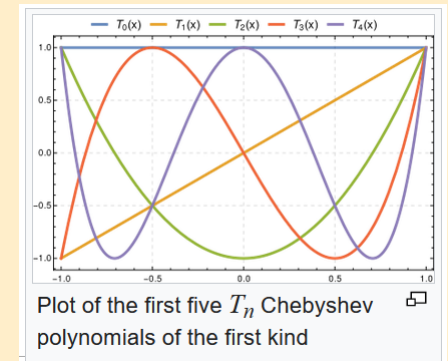
Solution is to approximate  $g_\theta(\Lambda)$  with the Chebyshev polynomials

Chebyshev series is related to a Fourier cosine series (period  $2L$ )

$$f(x) = \frac{c_0}{2} + \sum_{n=1}^{\infty} c_n \cos \frac{n\pi x}{L} \quad c_n = \frac{2}{L} \int_0^L f(x) \cos \frac{n\pi x}{L} dx, n \in \mathbb{N}_0.$$

$$\begin{aligned} T_0(x) &= 1 \\ T_1(x) &= x \\ T_{n+1}(x) &= 2x T_n(x) - T_{n-1}(x). \end{aligned}$$

$$\sum_{n=0}^{\infty} T_n(x) t^n = \frac{1 - tx}{1 - 2tx + t^2}.$$



Wikipedia

$$g_{\theta'}(\Lambda) \approx \sum_{k=0}^K \theta'_k T_k(\tilde{\Lambda})$$

$$\tilde{\Lambda} = \frac{2}{\lambda_{\max}} \Lambda - I_N$$

, where:

$\lambda_{\max}$  is the largest eigenvalue of the Laplacian

$\theta'$  is the Chebyshev coefficients

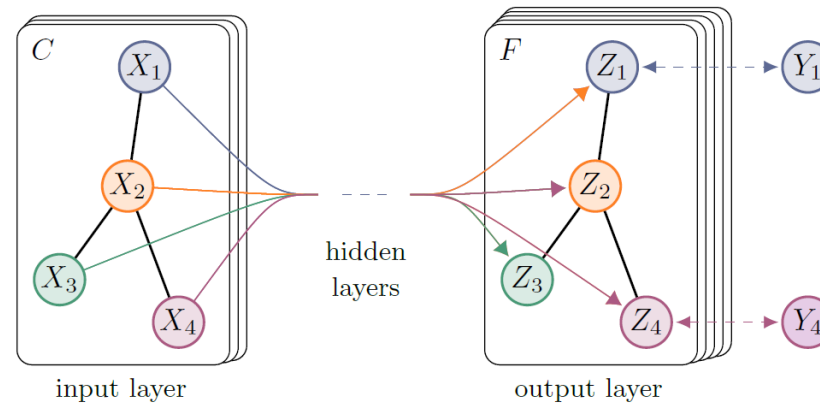
Chebyshev polynomial is recursively defined by

$$T_k(x) = 2xT_{k-1}(x) - T_{k-2}(x), \text{ where } T_0(x) = 1, T_1(x) = x$$

$$g_{\theta'} \star x \approx \sum_{k=0}^K \theta'_k T_k(\tilde{L})x$$

# Forward computation in the GNN [Kipf & Welling 2017]

C input channels  
F maps in the output layer  
 $Y_i$  label of node  $X_i$



(a) Graph Convolutional Network



(b) Hidden layer activations

$$Z = f(X, A) = \text{softmax}\left(\hat{A} \text{ReLU}\left(\hat{A}XW^{(0)}\right)W^{(1)}\right).$$

$$\text{softmax}(x_i) = \frac{\exp(x_i)}{\sum_i \exp(x_i)}$$

$$\mathcal{L} = - \sum_{l \in \mathcal{Y}_L} \sum_{f=1}^F Y_{lf} \ln Z_{lf},$$

Loss function is the cross-entropy, where  $Y_{lf}$  are the nodes with labels and

Kipf, Thomas N., and Max Welling. (2017) "Semi-supervised classification with graph convolutional networks." ICRL, *arXiv preprint arXiv:1609.02907*

Code to reproduce: <https://github.com/tkipf/gcn>

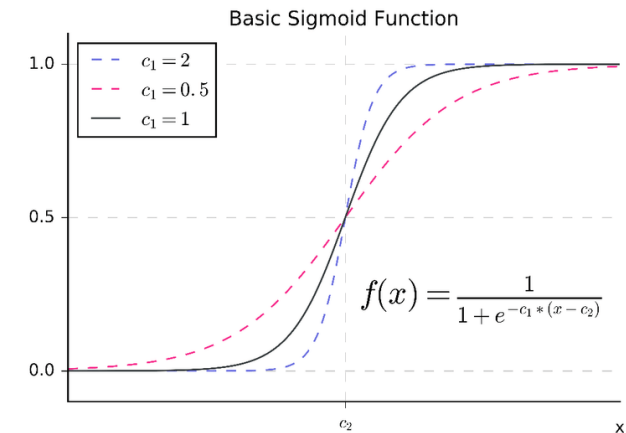
# Softmax [1]

$$h_{\theta}(x) = \frac{1}{1 + \exp(-\theta^{\top} x)}$$

where,  $x \in \mathbb{R}$  features and  $\theta$  parameters to optimize

$$P(y^{(i)} = k | x^{(i)}; \theta) = \frac{\exp(\theta^{(k)\top} x^{(i)})}{\sum_{j=1}^K \exp(\theta^{(j)\top} x^{(i)})}$$

where,  $y$  are the classes of  $x$



Cost function  $J$

$$J(\theta) = - \left[ \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) \right]$$

Gradient of  $J$

$$\nabla_{\theta^{(k)}} J(\theta) = - \sum_{i=1}^m \left[ x^{(i)} \left( \mathbf{1}\{y^{(i)} = k\} - P(y^{(i)} = k | x^{(i)}; \theta) \right) \right]$$

END