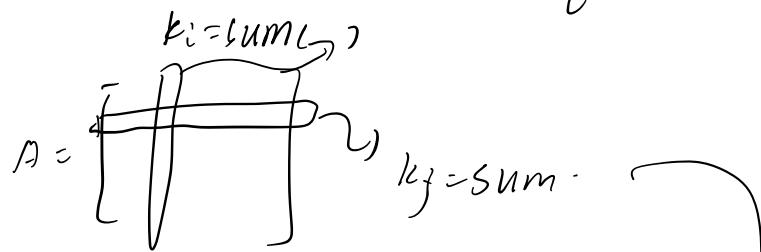


L1. Intro.

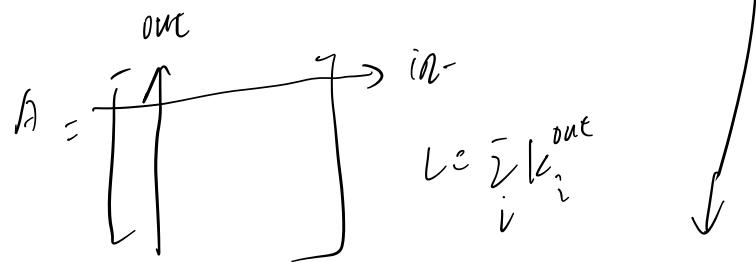
$$L = \frac{1}{2} \sum_i k_i$$

Adj Matrix :

• undirect:



• direct:

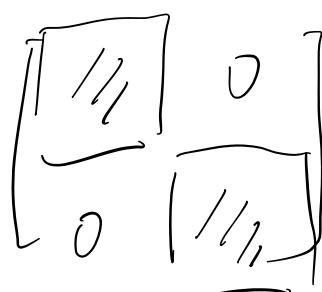


$$\text{平均度: } \bar{k} = \frac{2E}{N}$$

Weighted: $A = \begin{bmatrix} 0.1 & \cdots \\ 0.2 & \ddots \\ \vdots & \end{bmatrix}$

$$\bar{E} = \frac{1}{2} \cdot \sum_{ij}^N p_{ij} \quad \bar{k} = 2 \cdot \frac{\bar{E}}{N}$$

Disconnected:

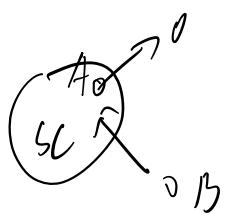


Connectivity:

• Strongly connected: 有且仅有2条路径 (A → B, B → A)

• Weakly ~ : ~ (不考虑方向)

SCC_S :



A: Out-component.

B: In-component

L2. Feature learning -

- ML need features.
- Goal: learn Structure - position of a node.

Traditional Methods: $\begin{cases} 1. \text{Degree} \\ 2. \text{centrality} \\ 3. \text{triangles} \end{cases}$

Node level:

- Degree : cons: can't capture importance of nodes.

- Centrality :

- Degree : $C_d(v_i) = d_i$

$$\begin{cases} C_d^{\text{norm}}(v_i) = \frac{d_i}{n-1} \\ C_d^{\text{max}}(v_i) = \frac{d_i}{\max_j(d_j)} \\ C_d^{\text{sum}}(v_i) = \frac{d_i}{2|E|} \end{cases}$$

- Betweenness :

$$C_B(v_i) = \sum_{j \in L} \frac{g_{jik}(i)}{g_{jk}} \quad \rightarrow \text{paths } j-k \text{ pass } i$$

$$g_{jk} \quad \rightarrow \text{all paths in } j-k$$

$$C_B'(v_i) = \frac{C_B(v_i)}{(n-1)(n-2)}$$

- Cluseness :

$$C_C(v_i) = \frac{n-1}{\sum_{j=1}^n d(i,j)}$$

Comparison:

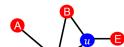
	Low Degree	Low Closeness	Low Betweenness
High Degree		Node is embedded in a community that is far from the rest of the network	Node's connections are redundant - communication bypasses the node
High Closeness	Key node connected to important/active alters		Probably multiple paths in the network, node is near many people, but so are many others
High Betweenness	Node's few ties are crucial for network flow	Very rare! Node monopolizes the ties from a small number of people to many others.	

- Graphlets:

- Def: **Induced subgraph** is another graph, formed from a subset of vertices and *all* of the edges connecting the vertices in that subset.

Induced subgraph:

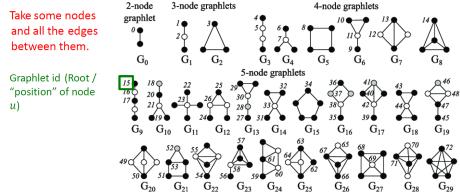
Not induced subgraph:



- Graphlets: Rooted connected induced non-isomorphic subgraphs:

There are 73 different graphlets on up to 5 nodes

Take some nodes and all the edges between them.

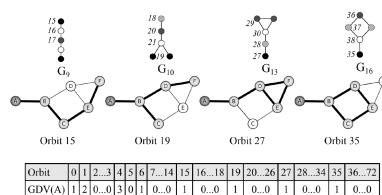


Note: Here is still on homogeneous graphs. Different colours distinguish different orbits and positions.

Homogeneous			Heterogenous		
Pattern	Graph	Count	Pattern	Graph	Count
		0			0
		12			4
		24			6
					1
					2

Graphlet Degree Vector

- Computation of the graphlet degree vector (GDV) of node A in the friendship network
 - GDV provides a measure of a node's local network topology
 - Comparing vectors of two nodes provides a more detailed measure of local topological similarity than node degrees



Orbit	0	1	2..3	4	5	6	7..14	15	16..18	19	20..26	27	28..34	35	36..72
GDV(A)	1	2	0..0	3	0	1	0..0	1	0..0	1	0..0	1	0..0	1	0..0

Link level: edge score : $C(x, y)$

1. Distance: $0 - v - 0$ \rightarrow $\frac{1}{2}$ \rightarrow $S_{AB} = 2$

2. Local Neighbourhood overlap:

$$\text{Jaccard} = \frac{|N(A) \cap N(B)|}{|N(A) \cup N(B)|}$$

13 PageRank.

$$\cdot r = Mr$$

$$\cdot \text{计算: } \hat{r}^k = [1/n, \dots, 1/n]^T$$

$$r^{k+1} = Mr^k$$

$$\text{直至: } \|r^{k+1} - r^k\|_1 < \varepsilon.$$

问题1: spider trap

$\rightarrow \text{②}\textcircled{2}$, 会导致 P_A 最高, 其它为 0

Solution:

$$r^{k+1} = \beta \cdot Mr^k + (1-\beta) \begin{bmatrix} 1/n \\ \vdots \\ 1/n \end{bmatrix}$$

(网页会随机跳到其它网页)

问题2: Dead End

$\rightarrow \text{②}$, 会导致 rank 全是 0. (Non-stochastic!)

$$\text{Solution: } \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \rightarrow \begin{bmatrix} 0 \\ \frac{1}{n} \\ \vdots \\ \frac{1}{n} \end{bmatrix}$$

问题3: 用稀疏矩阵. M 的空间太大. $M \approx 10N$ entries.

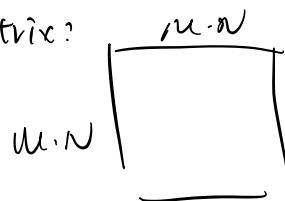
Solution: 利用 adj list, encode M .

$$\text{公式: } r \leftarrow \beta Mr + (1-\beta) \begin{bmatrix} 1/n \\ \vdots \\ 1/n \end{bmatrix}$$

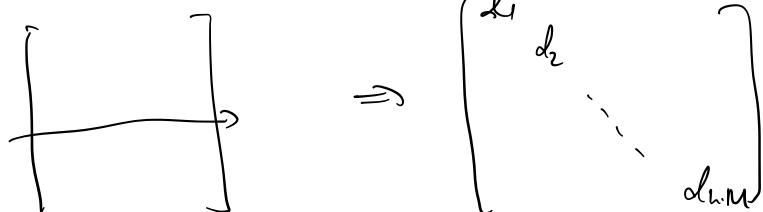
24. Graph Laplacian.

$$\text{Similarity: } w_{ij} = e^{-\frac{\|x_{ui} - x_{uj}\|_2^2}{\sigma_x^2}}$$

- Affinity Matrix?

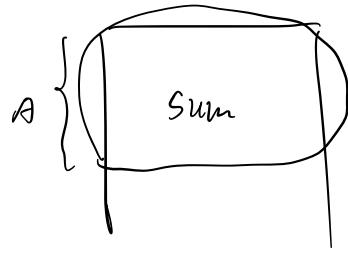


(1) Degree:



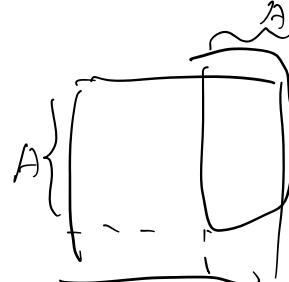
• Volume of set:

$$vol(A) = \sum_{i \in A} d_i$$



• Cuts:

$$cut(A, \bar{A}) = \sum_{i \in A, j \in \bar{A}} w_{ij}$$



- Laplace: $L = D - W$

• Spectral Clustering:

- Min Cut:

$$\text{MinCut}(A, B) = \min_{A, B} \sum_{u \in A, v \in B} w(u, v)$$

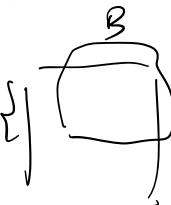
(intuition: 两团间连接尽可能少)

-Normalized Cut:

$$N\text{Cut}(A, B) = \text{Cut}(A, B) \left(\frac{1}{\text{Vol}(A)} + \frac{1}{\text{Vol}(B)} \right)$$

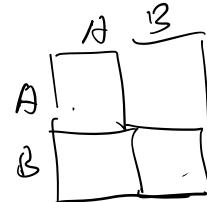
去掉连接到 outlier 的边

$$\text{association: } \text{assoc}(A, B) = \sum_{u \in A, v \in B} w(u, v)$$



$$N\text{assoc}(A, B) = \frac{\text{assoc}(A, A)}{\text{assoc}(A, B)} + \frac{\text{assoc}(B, B)}{\text{assoc}(A, B)}$$

$$\Rightarrow N\text{cut}(A, B) = 2 - N\text{assoc}(A, B)$$



Problem: 2-clustering.

Goal: find A, B . $\min_{A, B} N\text{cut}(A, B)$

NP-hard
↑

Solution: $\min_x N\text{cut}(x) = \min_y \frac{y^T(D-W)y}{y^T D y}$, subject to $y^T D y = 1$

$$\Leftrightarrow \min_y y^T(D-W)y \text{ subject to } y^T D y = 1$$

$$\Leftrightarrow \min_y y^T(D-W)y - \lambda(y^T D y - 1) = 0$$

$$\sum_i y_i D_{ii} = 0$$

$$(D-W)y = \lambda D y$$

$$\Leftrightarrow (D-W)y - \lambda D y = 0$$

$$\Rightarrow D^T(D-W)y = \lambda y$$

$$\Leftrightarrow (I - D^T W)y = \lambda y$$

\Rightarrow eigenvector.
(second smallest)

Avg association : $\max_{A,B} \left\{ \frac{\text{assoc}(A,A)}{|A|} + \frac{\text{assoc}(B,B)}{|B|} \right\}$



Avg cut : $\min_{A,B} \left\{ \frac{\text{Cut}(A,B)}{\sqrt{|A|}} + \frac{\text{Cut}(A,B)}{\sqrt{|B|}} \right\}$

Laplacian : $\mathcal{L} = D - W$.

$$\mathcal{L}_{\text{sym}} = D^{-\frac{1}{2}} (D - W) D^{-\frac{1}{2}}$$

$$\mathcal{L}_{\text{rw}} = D^{-1} (D - W)$$

$$\mathcal{L}_{\text{sys}} = D^{\frac{1}{2}} \mathcal{L}_{\text{rw}} D^{-\frac{1}{2}}$$

$y^T (D - W) y = \sum_{i,j} w_{ij} (y_i - y_j)^2$ — smoothness of G

$$y^T D^{-\frac{1}{2}} (D - W) D^{-\frac{1}{2}} y = \sum_{i,j} w_{ij} \left(\frac{y_i}{\sqrt{D_{ii}}} - \frac{y_j}{\sqrt{D_{jj}}} \right)^2$$

$(Dy)_i = \sum_{j \in N(i)} w_{ij} (y_i - y_j)$

\mathcal{L}^{-1} as operator:

$(D^{-1}W)^k$ k -th hop

$\sum_{k=0}^{\infty} (D^{-1}W)^k = I - D^{-1}W$

It's like this

$$(Ay)_i = \sum_j A_{ij} y_j$$

$$\Rightarrow y^T Ay = \left(\sum_i A_{ii} y_i \right) \cdot y_i \\ = \sum_{i,j} A_{ij} y_i y_j$$

$$y^T Ay = \sum_{i,j} A_{ij} y_i y_j$$

$$\frac{\partial y^T Ay}{\partial y_k} = \frac{\partial}{\partial y_k} \sum_{i,j} A_{ij} y_i y_j$$

$$= \sum_{j \neq k} A_{ij} y_j + 2A_{kk} y_k + \sum_{i \neq k} A_{ij} y_i \\ = \sum_j A_{kj} y_j + \sum_i A_{ik} y_i$$

$$\therefore \frac{\partial y^T Ay}{\partial y} = y^T D + y^T D^T = \underline{y^T (D+D^T)}$$

$$\frac{\partial y^T (D-W)y}{\partial y} = y^T (D-W + D^T - W^T)$$

$$= 2 y^T (D-W)$$

$$\Rightarrow (D-W)y - \lambda D y = v$$

$$y^T (D-W)y$$

$$(D-W)y_i = \sum_{j,i} (D-W)_{ij} y_j$$

$$\Rightarrow y^T (D-W)y = \left(\sum_{i,j} (D-W)_{ij} y_j \right) y_i$$

$$= \sum_{i,j} D_{ij} y_i y_j - \sum_{i,j} W_{ij} y_i y_j$$

$$= \sum_i D_{ii} y_i^2 - \sum_{i,j} W_{ij} y_i y_j$$

$$= \sum_i \sum_j W_{ij} y_i^2 - \sum_{i,j} W_{ij} y_i y_j + \sum_i \sum_j W_{ij}$$

$$= \sum_i \sum_j W_{ij} (y_i - y_j)^2$$

Cuts Comparison :

	Finding clumps	Finding splits	
Discrete formulation	Average association $\frac{\text{asso}(A,A)}{ A } + \frac{\text{asso}(B,B)}{ B }$	Normalized Cut $\frac{\text{cut}(A,B)}{\text{asso}(A,V)} + \frac{\text{cut}(A,B)}{\text{asso}(B,V)}$ or $2 - \left(\frac{\text{asso}(A,A)}{\text{asso}(A,V)} + \frac{\text{asso}(B,B)}{\text{asso}(B,V)} \right)$	Average cut $\frac{\text{cut}(A,B)}{ A } + \frac{\text{cut}(A,B)}{ B }$
Continuous solution	$Wx = \bar{\lambda}x$	$(D-W)x = \bar{\lambda}Dx$ or $Wx = (1-\bar{\lambda})Dx$	$(D-W)x = \bar{\lambda}x$

$$\begin{aligned}
 & f^T D^{-1} (D - W) D^{-1} f \\
 &= f^T f - f^T D^{-1} W D^{-1} f \\
 &= \sum_i^N f_i^2 - \left[\frac{f_i}{\sqrt{D_{ii}}} \right] \cdot \left[\sum_j^N \frac{f_j}{\sqrt{D_{jj}}} \cdot w_{ij} \right] \rightarrow i\text{-th row}
 \end{aligned}$$

$$\begin{aligned}
 &= \sum_i^N f_i^2 - \sum_i^N \frac{f_i}{\sqrt{D_{ii}}} \cdot \sum_j^N \frac{f_j}{\sqrt{D_{jj}}} w_{ij} \\
 &= \sum_i^N \sum_j^N f_i^2 - \frac{f_i \cdot f_j}{\sqrt{D_{ii}} \cdot \sqrt{D_{jj}}} \cdot w_{ij} \\
 &= \frac{1}{2} \cdot \sum_i^N \sum_j^N w_{ij} \left(\frac{f_i^2}{D_{ii}} - \frac{2f_i \cdot f_j}{\sqrt{D_{ii}} \sqrt{D_{jj}}} + \frac{f_j^2}{D_{jj}} \right) \\
 &= \frac{1}{2} \sum_{i,j} w_{ij} \left(\frac{f_i}{\sqrt{D_{ii}}} - \frac{f_j}{\sqrt{D_{jj}}} \right)^2
 \end{aligned}$$

- Personalized PageRank:

$$P^{t+1} = (1-\beta)q + \beta W P^t \quad (\text{不同的是, 它有一定概率回到原点})$$

• 用于 semi-supervised learning. (回到有label的地方)

• e.g.: (2 classes)

$$\text{Formula: } f^{t+1} = (1-\lambda)y + \lambda \underbrace{s f^t}_{D^{-1} W D^{-1}}$$

$$\Rightarrow f^{t+1} = (1-\lambda)y + \lambda s((1-\lambda)y + \lambda s f^{t+1})$$

$$= (1-\alpha)y \cdot \sum_{i=0}^{\infty} (\alpha s)^i + (\alpha s)^t y^v$$

$$\lim_{t \rightarrow \infty} \sum_i (\alpha s)^i = \frac{(1-\alpha s)^{t+1}}{1-\alpha s} = (1-\alpha s)^{-1}$$

$$\therefore f^* = \lim_{t \rightarrow \infty} f^t = (1-\alpha) \cdot y / (1-\alpha s)^{-1}$$

$$\begin{aligned} \text{In fact: } L(f) &= \frac{1}{2} \left(\sum_{ij} w_{ij} \left| \frac{f_i}{D_{ii}} - \frac{f_j}{D_{jj}} \right|^2 + \mu \sum_i \|f_i - y_i\|^2 \right) \\ &= \frac{1}{2} \left(f^T D^{-1} (D - w) D^{-1} f + \mu \cdot (f - y)^T (f - y) \right) \end{aligned}$$

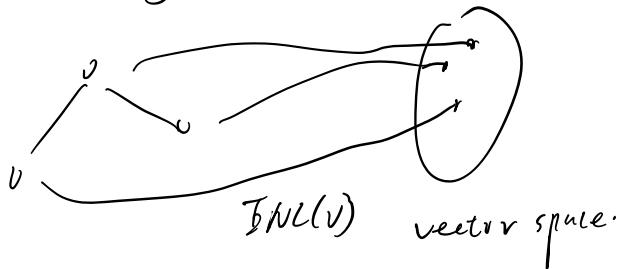
$$\frac{\partial L}{\partial f} = 0 \iff \underbrace{D^{-1} (D - w) D^{-1} f + \mu \cdot (f - y)}_A = 0$$

$$\iff f^* = \mu y \cdot (D^{-1} (D - w) D^{-1} + \mu \cdot I)^{-1}$$

(这即是 GCN 在前半部分 semi-supervised learning 的普遍做法)

25. Node / Graph Embedding.

A. Node Embedding.



- $\text{sim}(z_i, z_j) \approx z_i^T z_j$
- learn structure feature, embedding for downstream task.
- 2种寻找 nodes 之间关系的方法：

① Adj-based similarity:

$$- L = \sum_{(u,v) \in V \times V} \|z_u^T z_v - A_{u,v}\|^2$$

- 用 SGD 更新。

- Cons: 1. $O(V^2)$ time. (改由 sparse 图, 只对 non-zero
元素)
2. 只考虑 local, direct 邻接

② Random Walk:

- why? 1. Expressive: 考虑了 local, high-order 的信息。
2. 只考虑 Rand Walk 上的 node

• $N_R(u)$ — set of neighbour nodes. 12 - strategy

• objective: $\max_f \sum_{u \in V} \log P(N_R(u) | z_u)$

$$\Leftrightarrow L = - \sum_{u \in V} \sum_{v \in N(u)} \log P(v | z_u)$$

where $P(v | z_u) = \frac{\exp(z_v^T z_u)}{\sum_{k \in V} \exp(z_k^T z_u)}$

$$\Rightarrow L = - \sum_{u \in V} \sum_{v \in N(u)} \log \frac{\exp(z_v^T z_u)}{\sum_{k \in V} \exp(z_k^T z_u)}$$

$O(|V|^2)!$

• Word2Vec 2 step model: v_1, v_2, v_3, v_4, v_5

1. CBOW: [背景词] \rightarrow [中心词]

$$\begin{matrix} p^{v_2} & \xrightarrow{\quad z_2 \quad} & z_2 & \xrightarrow{\hat{z} = \frac{z_2 + z_4}{2}} & \hat{y} & \sim y_3 \\ p^{v_4} & \xrightarrow{\boxed{2 \times n}} & z_4 & \xrightarrow{\quad z_4 \quad} & & \end{matrix}$$

2. SK: [背景词] \rightarrow [中心词]

$$p^{v_3} \xrightarrow{z_{v \times n}} z_3 \xrightarrow{z_{n \times v}} \hat{y} \sim y_2 \text{ or } y_4$$

• 如何解决 $O(NV^2)$?

• Approximation:

$$\log \frac{\exp(z_v^\top z_u)}{\sum_{n \in V} (z_n^\top z_u)} \approx \log \delta(z_u^\top z_v) - \overbrace{\frac{1}{k} \log \delta(z_u^\top z_k)}^{\text{negative sample}}$$

• $k \mathcal{P}$, robust \mathcal{P}

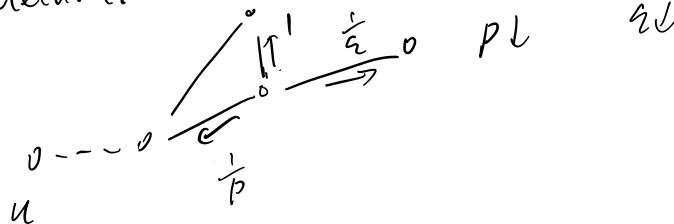
• k nodes are not on the walks!

• Better strategy: Biased Walks.

• why? trade off between local / global

\downarrow
BFS
 \downarrow
DFS

• Method:

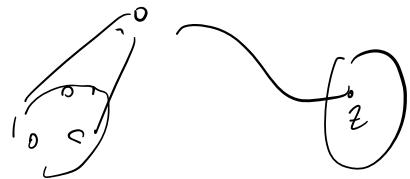


$$\pi_{C_i | C_{i1}, C_{i2}} = \begin{cases} 1 & , d(C_{i2}, C_i) = 1 \\ \alpha p & , \sim = 2 \\ \alpha q & , \sim = 0 \end{cases}$$

2. Graph Embedding:

• sum / avg of node embeddy $z_g = \sum_{v \in G} z_v$

\hat{v}_n : virtual node, and run node embedding



\hat{v}_n : Anonymous Walk Embedding.

- 1. do m times random walk with len L
- 2. record w_i counts
- 3. normalize w vector. w is G embedding.

$$w: \begin{array}{c} \boxed{\text{1} \text{ 1} \text{ 1} \text{ 1} \text{ 1}} \\ \diagup \quad \diagdown \quad \diagup \quad \diagdown \quad \diagup \\ \text{111} \quad \text{112} \quad \text{121} \quad \text{122} \quad \text{123} \end{array}$$

- How many walks?

$$m = \left\lceil \frac{2}{L^2} (\lg(\nu^{n-2}) - \lg(\nu)) \right\rceil$$

↑ theoretically existed types of
asyn. walks of len L

↓

$P(\# \text{errors} \leq \epsilon) < \delta$

2b GNN.

- Limitation of shallow Embedding:
 1. O(V) params
 2. can't generate new embedding for new nodes.
 3. no data feature.
- Goal: find $f: V \rightarrow \mathbb{R}^d$

- Naive: $MLP([A_i, f_i])$
 1. ~~O(V)~~: $O(V)$ parameters
 2. not invariant to node ordering
- GNN model should obey Permutation Equivariant/Invariant;

1. Permutation Equivariant:

同 - $\hat{\text{节点}}$, 不同 $\hat{\text{order}}$ $(A, X), (A', X')$, $f: G \rightarrow \mathbb{R}^d$:

对于同一个 node n : $f(A, X)_n = f(A', X')_n$

$\Leftrightarrow f$ is permutation equivariant.

2. Permutation Invariant:

同 - $\hat{\text{节点}}$, 不同 $\hat{\text{order}}$ $(A, X), (A', X')$, $f: G \rightarrow \mathbb{R}$:

$$f(A, X) = f(A', X')$$

$\Leftrightarrow f$ is Permutation Invariant.

• 2 Models:

1. Basic GNN:

• Math:

$$h_v^0 = \pi_v$$

$$h_v^k = \sigma \left(W_k \sum_{u \in N(v)} \frac{h_u^{k-1}}{|N(v)|} + B_k \cdot h_v^{k-1} \right), \quad \forall k > 0$$

• Matrix Formulation:

$$H^{(k+1)} = \sigma \left(\underbrace{D^{-1} A H^{(k)}}_T W_k^\top + H^k B_k^\top \right)$$

↓
projection

AH : aggregation.

• Train:

1. unsupervised: 利用结构，node similarity

• DeepWalk, Adj sim

2. supervised: $f: \mathbb{R}^n \rightarrow \mathbb{R}$. classification.

~ Pros: 1. share paras
2. inductive

2. GraphSage:

$$h_v^{(k)} = \sigma \left([W_k \cdot \text{Avg}(h_u^{(k-1)}, \text{but}(v)), B_k \cdot h_v^{(k-1)}] \right)$$

- 变化:
 1. Add(L) generalization
 2. sum \rightarrow concatenation.
- 3. Add AGG:
 1. Mean \rightarrow elementwise mean/max
 2. Pool: $y \left(\{R^{h_n^u}, \text{argmax} \} \right)$
 3. LSTM: $LSTM \left([h_n^{k1}, \underbrace{\text{the } \pi(N^V)}_{\text{生成一个随机序列}}] \right)$

L7. GNN 2.

Graph Signal: $f: V \rightarrow \mathbb{R}^{N \times d}$

Laplace Operator:

$$1. h = Lf \quad h[i] = \sum_{j \in N(V_i)} (f[i] - f[j])$$

$$2. f^T L f = \frac{1}{2} \sum_{i,j=1}^N \alpha_{ij} (f_i - f_j)^2$$

$$L = U \Lambda U^T$$

$$\hookrightarrow \begin{bmatrix} \lambda_0 & & \\ & \ddots & \\ & & \lambda_{N-1} \end{bmatrix}, \text{ where } 0 = \lambda_0 \leq \lambda_1 \leq \dots \leq \lambda_{N-1}.$$

$L^{-1} = 0 \cdot I$ (trivial solution)

$u_0 \dots u_m$ are basis signals, $\lambda_0 \dots \lambda_{N-1}$ are corresponding frequency.

Graph Fourier Transform: smoothness: $Lu = \sum_{i,j} \alpha_{ij} (u_i - u_j)$

$$L = U \Lambda U^T \quad (Lu = \lambda u) \Rightarrow |u| = 1$$

$$GFT : \hat{f} = U^T f$$

$$\text{Inverse GFT} : f = U \hat{f}$$

GNN Framework:

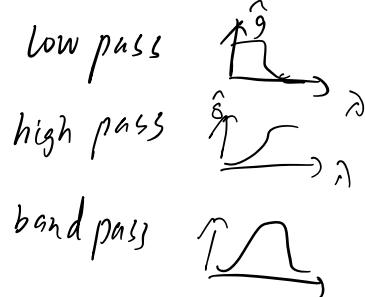
Filter: $A \in \mathbb{R}^{n \times n}, X \in \mathbb{R}^{n \times d} \Rightarrow A \in \mathbb{R}^{n \times n}, X \in \mathbb{R}^{n \times d_{\text{new}}}$

Pooling: $A \in \mathbb{R}^{n \times n}, X \in \mathbb{R}^{n \times d} \Rightarrow A \in \mathbb{R}^{n_p \times n_p}, X \in \mathbb{R}^{n_p \times d_{\text{new}}}$

Spectral Filtering:

$$\cdot f \xrightarrow{\text{GFT}} U^T f \xrightarrow{\hat{g}(\lambda)} g(\lambda) U^T f \xrightarrow{\text{IGFT}} U \hat{g}(\lambda) U^T f$$

$$\cdot \hat{g}(\lambda) = \begin{bmatrix} \hat{g}(\lambda_0) & & \\ & \ddots & \\ & & \hat{g}(\lambda_{n-1}) \end{bmatrix}$$



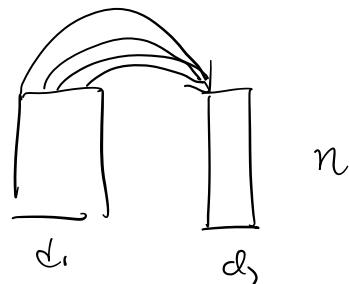
$$\cdot f \leftarrow \underbrace{U \hat{g}(\lambda) U^T}_{\hat{g}(z)} f$$

• we learn $g(\lambda)$ from data?

Multi-channel signals:

$$\tilde{f}_{in} \in \mathbb{R}^{n \times d_1} \rightarrow \tilde{f}_{out} \in \mathbb{R}^{n \times d_2}$$

$$\tilde{f}_{out}[:, i] = \sum_{j=1}^{d_1} \hat{g}_{i,j}(L) \cdot \tilde{f}_{in}[:, j], \quad i = 1, \dots, d_2$$



∴ we have $d_2 \times d_1$ filters.

$$\cdot g(\lambda) = \begin{bmatrix} \theta_1 & \dots & \\ \vdots & \ddots & \\ \theta_n & & \end{bmatrix} \Rightarrow d_2 \times d_1 \times N \text{ parameters.}$$

- Polynomial Parameterized:

$$\hat{g}(n) = \begin{bmatrix} \sum_{k=0}^K b_k \lambda_1^k \\ \vdots \\ \sum_{k=0}^K b_k \lambda_N^k \end{bmatrix}$$

$$L^k = (U \Lambda U^T)^k = U \Lambda^k U^T \rightarrow \begin{bmatrix} \lambda_1^k \\ \vdots \\ \lambda_N^k \end{bmatrix}$$

- $d_1 \times d_1 \times k$ params.

$$U \hat{g}(n) U^T f = \sum_{k=0}^K b_k L^k f \quad (\text{no eigen-decomposition needed})$$

- Spatial View: $L^k f$: smoothness within k -hops.

- Chebyshev Polynomial:

$$- T_0(x) = 1, \quad T_1(x) = x, \quad T_{n+1}(x) = 2 \cdot T_n(x) - T_{n-1}(x)$$

- Moivre's formula: $(\cos x + i \sin x)^n = \cos nx + i \sin nx$
(proved by induction.)

- Orthogonal polynomial $\stackrel{\text{def}}{\Leftarrow}$

WM integrable - $W(x) > 0$, $\langle f_m, f_n \rangle = \int_a^b f_m(x) \cdot f_n(x) W(x) dx$

$$\langle f_m, f_n \rangle = 0 \quad \text{for } m \neq n.$$

• Gram-Schmidt orthonormalization:

$$\text{proj}_u(v) \stackrel{\text{def}}{=} \frac{\langle u, v \rangle}{\langle u, u \rangle} \cdot \hat{u}$$

$$u_1 = v_1$$

$$u_2 = v_2 - \text{proj}_{u_1} v_2$$

$$u_3 = v_3 - \text{proj}_{u_1} v_3 - \text{proj}_{u_2} v_3$$

$$u_k = v_k - \sum_{i=1}^{k-1} \text{proj}_{u_i} v_k$$

then normalize $\{u_i\}_{i=1}^k$

• Def by trigonometric function:

$$T_n(\cos \theta) = \cos(n\theta)$$

• Proof correctness of Chebyshev polynomials are orthogonal:

$$\begin{aligned} \langle T_n, T_m \rangle &= \int_{-1}^1 T_n(x) \cdot T_m(x) \underbrace{\frac{1}{\sqrt{1-x^2}} dx}_{w(x)} \\ &= \int_{-1}^1 \cos(n \arccos x) \cdot \cos(m \arccos x) \underbrace{\frac{1}{\sqrt{1-x^2}} dx}_{w(x)} \\ &= \int_{-\pi}^{\pi} \cos(n \arccos x) \cdot \cos(m \arccos x) dx \arccos x \\ &= \int_{-\pi}^{\pi} \cos(ny) \cos(my) dy \end{aligned}$$

$$= \int_0^0 \left(m_3((m-n)y) + m_3((m+n)y) \right) dy \\ = 0 \quad (m \neq n).$$

!~ orthogonal!

普通多项式，x值可能特别大/小，造成数值问题

• ChebNet: (stable under perturbation of coefficients)

$$\xrightarrow{x \in [-1, 1], y \in [-1, 1]}$$

• Chebyshev polynomial: $T_0 = 1, T_1 = x, T_{n+1}(x) = 2T_n(x) - T_{n-1}(x)$

$$g(x) = \theta_0 T_0(x) + \theta_1 \cdot T_1(x) + \dots$$

$\xrightarrow{\text{Initialization for Cheb poly, } x \in [-1, 1]}$

$$\hat{\Lambda} = \frac{2I}{\lambda_{\max}} - I, \quad \hat{g}(\hat{\Lambda}) = \sum_{k=0}^K \theta_k T_k(\hat{\Lambda}) \quad \xrightarrow{\text{recursive def}}$$

$$\therefore u \hat{g}(\hat{\Lambda}) u^T f = \sum_{k=0}^K \theta_k T_k(L) f, \text{ where } L = \frac{2I}{\lambda_{\max}} - I.$$

• GCN: set $K=1, \lambda_{\max}=2$ for ChebNet.

$$\hat{g}(\hat{\Lambda}) = \theta_0 + \theta_1 (\Lambda - I) \quad \xrightarrow{\text{generally we want GCN be a low-pass filter for us}}$$

$$\text{let } -\theta_1 = \theta_0 = \theta \quad (\text{further constraint})$$

$$\hat{g}(\hat{\Lambda}) = \theta(2I - \Lambda)$$

$$\begin{aligned} u \hat{g}(\hat{\Lambda}) u^T f &= \theta(2I - L)f = \theta(2I - D^{-\frac{1}{2}}(D-A)D^{-\frac{1}{2}})f \\ &= \theta(I + D^{-\frac{1}{2}}A D^{-\frac{1}{2}})f \end{aligned}$$

$$\text{Renormalization: } u \hat{g}(\hat{\Lambda}) u^T f = \theta(D^{-\frac{1}{2}}\tilde{A} D^{-\frac{1}{2}})f, \text{ where } \tilde{A} = A + I$$

• multi-channel GIN:

$$\text{In general: } F_{\text{out}}[:, i] = \sum_{j=1}^{d_1} g_{i,j}(L) \cdot F_{\text{in}}[i, j], i = 1 \dots d_2.$$

$$\text{GIN : } F_{\text{out}}[i, i] = \sum_{j=1}^{d_1} \theta_{i,j}(D^{-\frac{1}{2}} A D^{\frac{1}{2}}) F_{\text{in}}[:, j] \quad i = 1 \dots d_2$$

$$\text{Matrix form: } F_{\text{out}} = (D^{-\frac{1}{2}} A D^{\frac{1}{2}}) F_{\text{in}} \cdot \theta, \quad \theta \in \mathbb{R}^{d_1 \times d_2}.$$

• Spatial view for GIN:

$$\text{let } C = D^{-\frac{1}{2}} A D^{\frac{1}{2}}$$

$$F_{\text{out}} = C F_{\text{in}} \theta$$

we know that: $C[i, j]$ if $v_j \notin N(v_i) \cup \{v_i\}$.

$$\Rightarrow F_{\text{out}}[i, :] = \sum_{v_j \in N(v_i) \cup \{v_i\}} C[i, j] \cdot \underbrace{F_{\text{in}}[j, :]}_{\text{Aggregation}} \theta.$$

Feature transformation

• GAT:

$$\alpha_{ij} = \frac{\exp(\text{G}(a^T [wh_i, wh_j]))}{\sum_{k \in N_i} \exp(\text{G}(a^T [wh_i, wh_k]))}$$



G is Leaky ReLU.

MPNN: (edge embedding)

$$m_i^{(k)} = \sum_{v_j \in N(v_i)} M_k(h_i^{(k)}, h_j^{(k)}, e_{ij})$$

$h_i^{(k+1)} = U_k(h_i^{(k)}, m_i^{(k)})$, where M_k, U_k need to be designed.

Graph Pooling: $\begin{cases} gPool: \text{节点重要性} \\ \text{DiffPool:} \text{聚类} \\ \text{EigenPooling: 利用特征器 GFT, 得到更好的 Hidden} \end{cases}$

H_f

1. $gPool: v_i \rightarrow y_i, y_i = \frac{h_i^T p}{\|p\|}$ Importance Ranking

$$idx \leftarrow \text{rank}(y, np)$$

$$A_p \leftarrow A[idx, idx] \quad \text{new Adj Matrix}$$

$$H_p \leftarrow H[i, idx, :] \odot \sigma(y[idx]) \quad \text{new hidden}$$

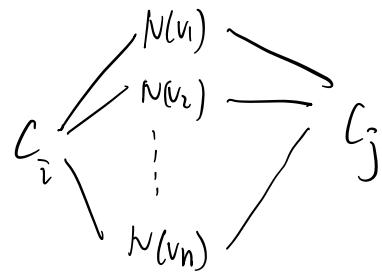
cons: 反映群聚性，使得 sample 和 node 相关性不强

2. $DiffPool:$

$$\begin{array}{ccc} H \in R^{n \times d} & \xrightarrow{w_1} & H_a \in R^{n \times np} \\ & \xrightarrow{w_2} & H_f \in R^{n \times d_{\text{new}}} \end{array}$$

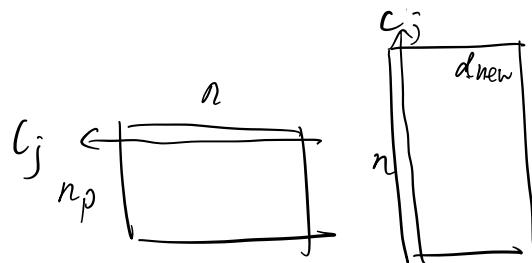
$$\begin{array}{c} \xrightarrow{} A_p = H_a^T A H_a \\ \xrightarrow{} H_p = H_a^T H_f \end{array}$$

$$\text{原理: } 1. \mathbf{A}p = \mathbf{H}\mathbf{a}^T \mathbf{A} \mathbf{H}\mathbf{a}$$



cluster \longleftrightarrow nodes & $N(v) \longleftrightarrow$ cluster.

$$2. \mathbf{H}_p = \mathbf{H}\mathbf{a}^T \mathbf{H}_f$$



Hidden feature of C_j is from $\underbrace{\text{weighted sum of } H_p}_{j}$

"weight" is prob of Node i & $N(i)$ in C_j

3. Eigen pooling: Step 1: learn $\mathbf{A}p$ like DiffPool

Step 2: use GFT to learn a better H_p .

subgraph $G_i \xrightarrow{\text{GFT}} \hat{f} \rightarrow \text{truncation some } \hat{f}_i$

X

- Eigen-decomposition is $O(N^3)$!!!

- If we want to develop a low-pass filter simply:

$$\hat{g}(\lambda) = \begin{bmatrix} \lambda_{\max} - \lambda_0 & = \lambda_{\max} \\ & \lambda_{\max} - \lambda_1 \\ & \ddots \\ & \lambda_{\max} - \lambda_N \end{bmatrix}$$

\leftarrow low f
 \leftarrow u/f
 \leftarrow high f

X

- Empirically Leaky ReLU > ReLU

· 18 GNN 3:

- **GNN Layer:**
 - Transformation + Aggregation
 - Classic GNN layers: GCN, GraphSAGE, GAT
- **Layer connectivity:**
 - Deciding number of layers
 - Skip connections
- **Graph Manipulation:**
 - Feature augmentation
 - Structure manipulation

1. GNN layer:

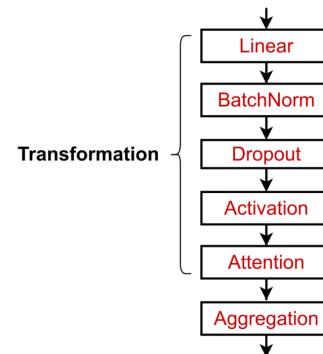
A Single GNN Layer

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



DNN tricks: (in a linear layer)

A suggested GNN Layer



2. More layers

- Add More layers \Rightarrow Over Smooth Problem. (Embedding are the same for different nodes)



- **Receptive field:** the set of nodes that determine the embedding of a node of interest

- We knew the embedding of a node is determined by its **receptive field**
 - If two nodes have highly-overlapped receptive fields, then their embeddings are highly similar
- Stack many GNN layers \rightarrow nodes will have highly-overlapped receptive fields \rightarrow node embeddings will be highly similar \rightarrow suffer from the over-smoothing problem

3. Designing GNNs: 求合理的层数。

- Step 1: Analyze the necessary receptive field to solve your problem. E.g., by computing the diameter of the graph
- Step 2: Set number of GNN layers L to be a bit more than the receptive field we like. Do not set L to be unnecessarily large!

4. Design Model (设计模型):

① Skip-connection Formula

A standard GCN layer

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

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

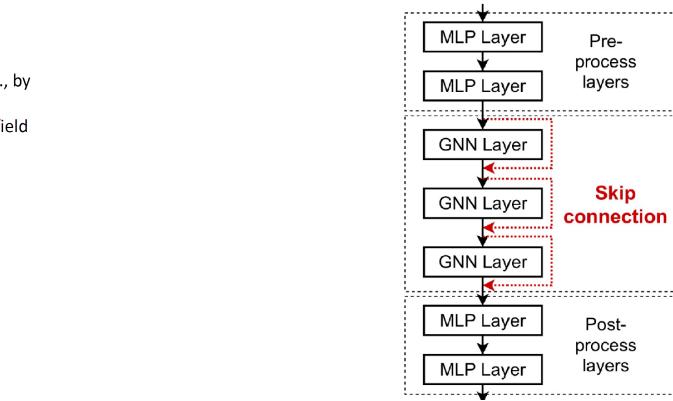
A GCN layer with skip connection

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

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

Other commonly used augmented features:

- PageRank
- Centrality



② 加强 GNN 的表达能力

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

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

③ GNN layer 设计. 单个 layer 加深

3. Augmentation. } Feature Structure

• Graph Feature manipulation

The input graph lacks features → feature augmentation

• Graph Structure manipulation

The graph is too sparse → Add virtual nodes / edges
The graph is too dense → Sample neighbors when doing message passing
The graph is too large → Sample subgraphs to compute embeddings

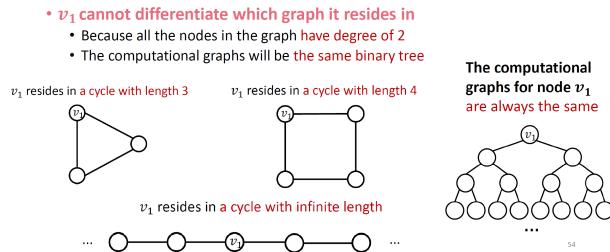
1. Feature Augment

	Constant node feature	One-hot node feature
Expressive power	Medium. All the nodes are identical, but GNN can still learn from the graph structure	High. Each node has a unique ID, so node-specific information can be stored
Inductive learning (Generalize to unseen nodes)	High. Simple to generalize to new nodes: we assign constant feature to them, then apply our GNN	Low. Cannot generalize to new nodes: new nodes introduce new IDs, GNN doesn't know how to embed unseen IDs
Computational cost	Low. Only 1 dimensional feature	High. $O(V)$ dimensional feature, cannot apply to large graphs
Use cases	Any graph, inductive settings (generalize to new nodes)	Small graph, transductive settings (no new nodes)

Why Feature Aug?

① Input hasn't any features, more expressive (系數)

② Some structures are hard to learn



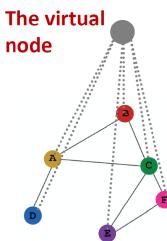
other Feature Aug Tricks: PageRank, Centrality ...

2) structure augmentation

I), Add virtual edges: K hops: $\tilde{A} = \sum K^k$

II), Add virtual nodes:

improve message passing in sparse graph.



Why we need these augmentation?

- Our assumption so far has been **Raw input graph = computational graph**

• Reasons for breaking this assumption

• Feature level:

- The input graph **lacks features** → feature augmentation

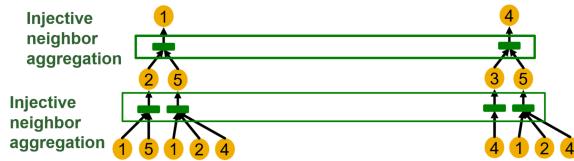
• Structure level:

- The graph is **too sparse** → inefficient message passing
- The graph is **too dense** → message passing is too costly
- The graph is **too large** → cannot fit the computational graph into a GPU

29 Graph Isomorphism Network.

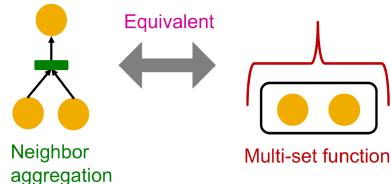
(研究的是：GNN对子结构的表达能力)

- More expressive: Injective neighbour aggregation:



Every step of aggregation is injective \Rightarrow GNN is injective \Rightarrow most expressive

- Multiset-function:



Injective function can distinguish multi-sets!

- Fail Examples:

- GCN** (mean-pool) [Kipf & Welling, ICLR 2017]
 - Uses element-wise mean pooling over neighboring node features

$$\text{Mean}(\{x_u\}_{u \in N(v)})$$

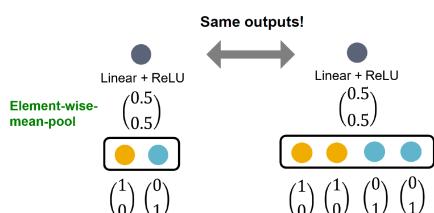
- GraphSAGE** (max-pool) [Hamilton et al., NeurIPS 2017]
 - Uses element-wise max pooling over neighboring node features

① GCN:

Theorem [Xu et al., ICLR 2019]

- GCN's aggregation function cannot distinguish different multi-sets with the same color proportion.

Failure case illustration

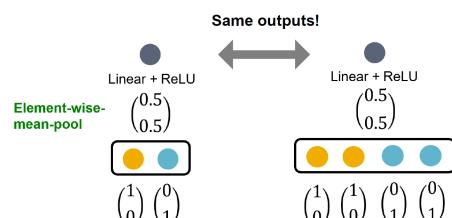


$$\text{Max}(\{x_u\}_{u \in N(v)})$$

② GraphSAGE (max pool)

- GraphSAGE's aggregation function cannot distinguish different multi-sets with the same set of distinct colors.

Failure case illustration



• Injective Multi-set function:

- Thm: Injective function can be expressed as: $\Phi\left(\sum_{x \in S} f(x)\right)$
- How to model $\Phi \cdot f$? MLP.

Thm (Universal Approximation Thm):

MLP ($n_{\text{hidden}} = 1$) with $\delta(\cdot)$ can approximate any continuous function.

$$\Rightarrow \text{Injective function} \Leftrightarrow \text{MLP}_{\Phi}\left(\sum_{x \in S} \text{MLP}_f(x)\right)$$

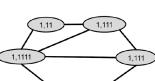
• GIN: $\text{MLP}_{\Phi}\left(\sum_{x \in S} \text{MLP}_f(x)\right)$

- Thm: GIN's aggregation is injective.
- \Rightarrow GIN is most expressive!

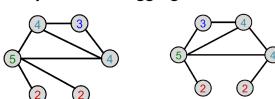
• WL graph kernel:

- Checking isomorphism is NP-hard!
- WL kernel can distinguish most of the real-world graphs.
- Color refinement: $C^{(k+1)}(v) = \text{HASH}(C^{(k)}(v), \{C(u)\}_{u \sim N(v)})$

• Aggregated colors:



• Injectively HASH the aggregated colors

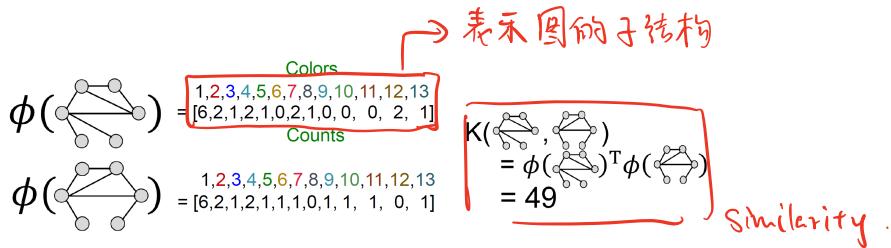


HASH table: Injective!

1,1	-->	2
1,11	-->	3
1,111	-->	4
1,1111	-->	5

HASH is injective!

- Iterate till: no difference between nodes in a Graph in iteration i..it.



- KL kernel can somehow show the similarity.
However $\phi(G_1) = \phi(G_2)$ doesn't imply isomorphism!
- $O(|E|)$ algorithm.
(since there are $2|E|$ multiplication for each iteration)
- GIN model: use NN to model the HASH function

$$\text{MLP}_\phi \left((1+\epsilon) \cdot \text{MLP}_f(c^k(v)) + \sum_{u \in N(v)} \text{MLP}_f(c^k(u)) \right)$$

- Since MLP can provide "one-hot" input for next layer:

$$\text{GINConv}(c^k(v), \{c^k(u)\}_{u \in N(v)}) = \text{MLP}_\phi \left((1+\epsilon) \cdot c^k(v) + \sum_{u \in N(v)} c^{(k+1)}(u) \right)$$

- Iteratively update node vectors by
 $c^{(k+1)}(v) = \text{GINConv}(\{c^{(k)}(v), \{c^{(k)}(u)\}_{u \in N(v)}\})$

- GIN can be understood as differentiable neural version of the WL graph Kernel:

	Update target	Update function
WL Graph Kernel	Node colors (one-hot)	HASH
GIN	Node embeddings (low-dim vectors)	GINConv

↳ capture fine-grained similarity of different nodes.

- Rank of expression: Sum > mean > max

L10 Heterogeneous Graphs, Knowledge Graphs.

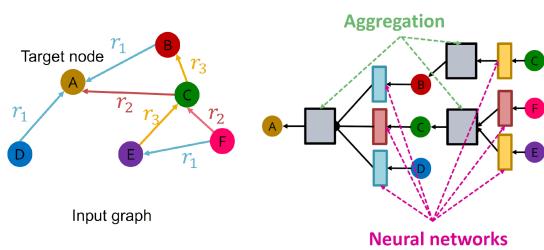
1. Heterogeneous Graphs

A heterogeneous graph is defined as

- $G = (V, E, R, T)$
- Nodes with node types $v_i \in V$
- Edges with relation types $(v_i, r, v_j) \in E$
- Node type $T(v_i)$
- Relation type $r \in R$

Method: Relational GLN

- Use different neural network weights for different relation types.



$$\cdot h_v^{(k+1)} = \text{ReLU} \left(\sum_{r \in R} \sum_{u \in N_v^r} \frac{1}{C_{v,r}} W_r^k h_u^{(k)} + W_0^k h_v^{(k)} \right)$$

$$\cdot \text{params: } \{W_r^1, \dots, W_r^L\}_{r \in R}$$

\Rightarrow too many params, overfitting

Solution:

① block diagonal Matrix:

- Key insight: make the weights sparse!
- Use block diagonal matrices for $W_r^{(l)}$

$$W_r = \begin{pmatrix} & & \\ & \text{green circles} & \\ & & \end{pmatrix}$$

Limitation: only nearby neurons/dimensions can interact through W

- If use B low-dimensional matrices, then # param reduces from $d^{(l+1)} \times d^{(l)}$ to $B \times \frac{d^{(l+1)}}{B} \times \frac{d^{(l)}}{B}$

② Basis learning:

$$W_r = \sum_{b=1}^B A_{rb} \cdot V_b \xrightarrow{\text{importance}} \text{basis matrix}$$

2. Knowledge Graph

• KG completion : given (h, r) , predict tails.

that is $\underset{\theta}{\operatorname{argmax}} \prod_{i,j,k \in N} p(X_{ijk}=1 | \theta, \theta_i)$, $X_{ijk}=1$ means $\rightarrow k$ is plausible.

• Methods:

• Path Rank Alg. (Random Walk)

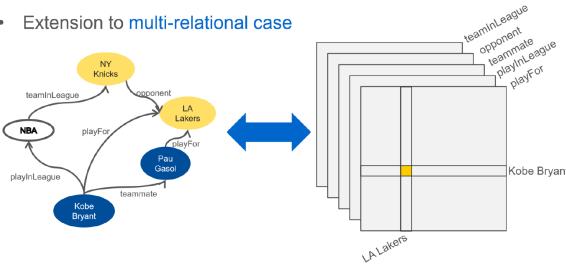
$$S(A, r, B) = \sum_{p \in \text{paths}(A \rightarrow B)} \theta_r^p \cdot \Pr(p)$$

$$\text{rel} = \underset{\text{relels}}{\operatorname{argmax}} S(A, r, B)$$

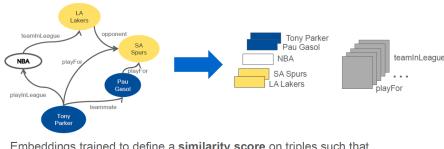
$$\text{obj} = \underset{\text{events}}{\operatorname{argmax}} S(A, R, e)$$

• Factorization:

- Extension to multi-relational case



- Related to Deep Learning methods
- Entities are vectors (low-dimensional sparse)
- Relation types are operators on these vectors



Embeddings trained to define a similarity score on triples such that

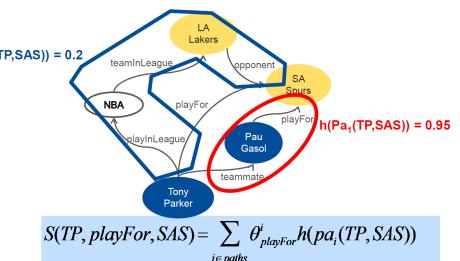
$$S(KB, playFor, LAL) > S(KB, playFor, NYK)$$

- KG Embedding : (h, r, t) .

Goal : $(h, r) \approx t$

Path Ranking Algorithm [Lao et al., 11]

- Random walks on the graph are used to sample paths
- Paths are weighted with probability of reaching target from source
- Paths are used as ranking experts in a scoring function



Relation Patterns

- **Symmetric (Antisymmetric) Relations:**

- **Example:**
 - Symmetric: Family, Roommate
 - Antisymmetric: Hyponym

$$r(h, t) \Rightarrow r(t, h) \quad (r(h, t) \Rightarrow \neg r(t, h)) \quad \forall h, t$$

- **Inverse Relations:**

- **Example :** (Advisor, Advisee)

$$r_2(h, t) \Rightarrow r_1(t, h)$$

- **Composition (Transitive) Relations:** $r_1(x, y) \wedge r_2(y, z) \Rightarrow r_3(x, z) \quad \forall x, y, z$

- **Example:** My mother's husband is my father.

- **1-to-N relations:**

$r(h, t_1), r(h, t_2), \dots, r(h, t_n)$ are all True.

- **Example:** r is "StudentsOf"

1. TransE: $h + r \approx t$

$$\text{Score} : f_r(h, t) = -d_r(h, t) = -\|h + r - t\|$$

$$\begin{aligned} \text{Loss function: } L(h, r, t) &= \max(0, \varphi + d_{r, \text{pos}} - d_{r, \text{neg}}) \\ &= \max(0, \varphi + f_{r, \text{neg}} - f_{r, \text{pos}}) \end{aligned}$$

Algorithm 1 Learning TransE

```

input Training set  $S = \{(h, \ell, t)\}$ , entities and rel. sets  $E$  and  $L$ , margin  $\gamma$ , embeddings dim.  $k$ .
1: initialize  $\ell \leftarrow \text{uniform}(-\frac{6}{\sqrt{k}}, \frac{6}{\sqrt{k}})$  for each  $\ell \in L$                                 Entities and relations are
2:  $\ell \leftarrow \ell / \|\ell\|$  for each  $\ell \in L$                                               initialized uniformly, and
3:  $e \leftarrow \text{uniform}(-\frac{6}{\sqrt{k}}, \frac{6}{\sqrt{k}})$  for each entity  $e \in E$                                               normalized
4: loop
5:  $e \leftarrow e / \|e\|$  for each entity  $e \in E$ 
6:  $S_{\text{batch}} \leftarrow \text{sample}(S, b)$  // sample a minibatch of size  $b$ 
7:  $T_{\text{batch}} \leftarrow \emptyset$  // initialize the set of pairs of triplets
8: for  $(h, \ell, t) \in S_{\text{batch}}$  do
9:    $(h', \ell, t') \leftarrow \text{sample}(S'_{(h, \ell, t)})$  // sample a corrupted triplet
10:   $T_{\text{batch}} \leftarrow T_{\text{batch}} \cup \{(h, \ell, t), (h', \ell, t')\}$  //  $d$  represents distance (negative of score)
11: end for
12: Update embeddings w.r.t.  $\sum_{((h, \ell, t), (h', \ell, t')) \in T_{\text{batch}}} \nabla [\gamma + d(h + \ell, t) - d(h' + \ell, t')] +$ 
13: end loop

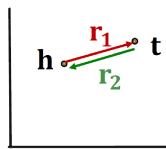
```

Contrastive loss: favors lower distance (or higher score) for valid triplets, high distance (or lower score) for corrupted ones

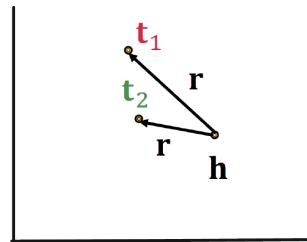
46

Inverse (\vee)

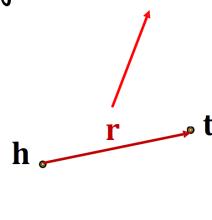
$h + r_2 = t$, we can set $r_1 = -r_2$



$1 - \tau_0 - N(x)$

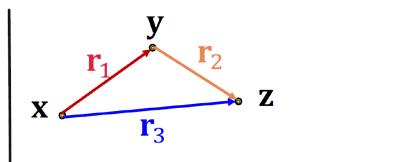


Symmetric (\bowtie)

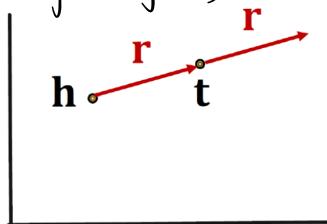


Composition (\wedge)

$$\bullet r_1 = r_2 + r_3$$



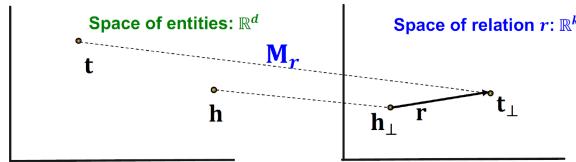
Asymmetry (\vee)



2. TransR: entity space $\mathbb{R}^d \xrightarrow{M_r r^k \text{ val}} \text{relation space } \mathbb{R}^k$

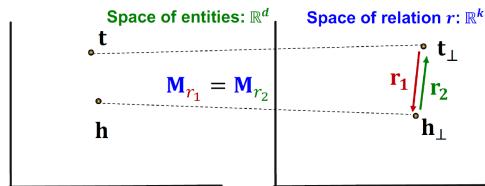
- $h_\perp = M_r h, t_\perp = M_r t$

- Score: $- \|h_\perp + r - t_\perp\|$



Inverse: ✓

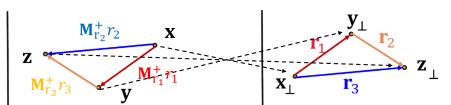
$$\begin{aligned} \mathbf{r}_2 &= -\mathbf{r}_1, \mathbf{M}_{r_1} = \mathbf{M}_{r_2} \\ \mathbf{M}_{r_1} \mathbf{t} + \mathbf{r}_1 &= \mathbf{M}_{r_1} \mathbf{h} \quad \mathbf{M}_{r_2} \mathbf{h} + \mathbf{r}_2 = \mathbf{M}_{r_2} \mathbf{t} \end{aligned}$$



Composition ✓

- TransR can model composition relations ✓

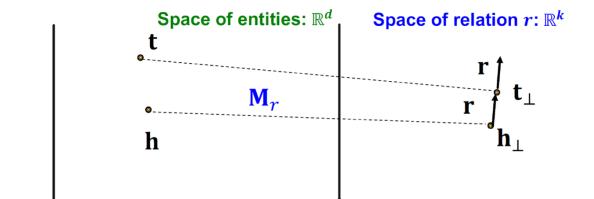
$M_{r_i}^+ = M_{r_i}^T (M_{r_i} M_{r_i}^T)^{-1}$ is the Moore-Penrose inverse (pseudoinverse) of M_{r_i} , where $M_{r_i} M_{r_i}^T$ is invertible (full rank); so that $M_{r_i} M_{r_i}^+ = I$



Asymmetry: ✓

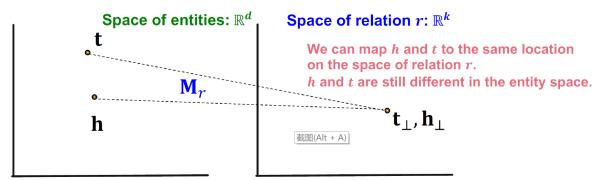
$$\mathbf{r} \neq 0, \mathbf{M}_r \mathbf{h} + \mathbf{r} = \mathbf{M}_r \mathbf{t},$$

Then $\mathbf{M}_r \mathbf{t} + \mathbf{r} \neq \mathbf{M}_r \mathbf{h}$ ✓



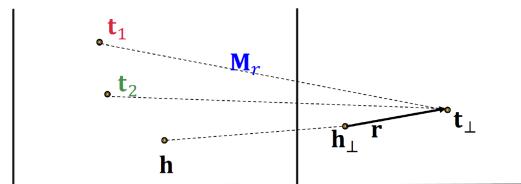
Symmetry: ✓

$$\mathbf{r} = 0, \mathbf{h}_\perp = \mathbf{M}_r \mathbf{h} = \mathbf{M}_r \mathbf{t} = \mathbf{t}_\perp \checkmark$$



1-t0-N: ✓

- We can learn \mathbf{M}_r so that $\mathbf{t}_\perp = \mathbf{M}_r \mathbf{t}_1 = \mathbf{M}_r \mathbf{t}_2$
- Note that \mathbf{t}_1 does not need to be equal to \mathbf{t}_2 !

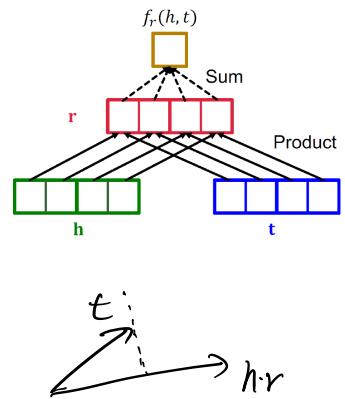


3. DistMult : (bilinear Modeling)

$$\cdot h, r, t \in \mathbb{R}^k$$

$$\cdot \text{Score: } f_r(h, t) = \langle h, r, t \rangle = \sum_i h_i \cdot r_i \cdot t_i$$

• Intuition: dot product of $h \cdot r$ and t .



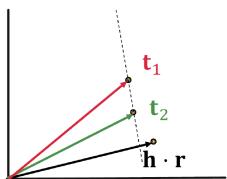
Symmetric: (✓)

DistMult can naturally model symmetric relations

$$f_r(h, t) = \langle h, r, t \rangle = \sum_i h_i \cdot r_i \cdot t_i = \langle t, r, h \rangle = f_r(t, h)$$

$$1\text{-to-N} \quad (\checkmark)$$

$$\langle h, r, t_1 \rangle = \langle h, r, t_2 \rangle$$



Inverse (✗)

DistMult cannot model inverse relations

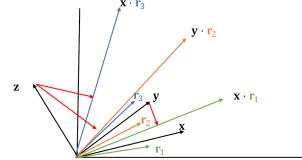
- If it does model inverse relations:

$$f_{r_2}(h, t) = \langle h, r_2, t \rangle = \langle t, r_1, h \rangle = f_{r_1}(t, h)$$

- This means $r_1 = r_2$

Composition (✗)

Intuition: **DistMult** defines a hyperplane for each (head, relation), the union of the hyperplanes induced by multi-hops of relations, e.g., (r_1, r_2) , cannot be expressed using a single hyperplane.



5 Embedding Alg:

Model	Sym.	Anisym.	Inv.	Compos.	1-to-N
TransE	✗	✓	✓	✓	✗
TransR	✓	✓	✓	✓	✓
DistMult	✓	✗	✗	✗	✓
ComplEx	✓	✓	✓	✗	✓
RotatE	✓	✓	✓	✓	✓

L10.11 Recommendation System.

1. Collaborative Filtering (CF)

- Input: user-item matrix
- Output: recommended items

💡:

- A popular similarity measure in user-based CF: **Pearson correlation**

$$a, b : \text{users}$$

$$r_{a,p} : \text{rating of user } a \text{ for item } p$$

$$\text{sim}(a, b) = \frac{\sum_{p \in P} (r_{a,p} - \bar{r}_a)(r_{b,p} - \bar{r}_b)}{\sqrt{\sum_{p \in P} (r_{a,p} - \bar{r}_a)^2} \sqrt{\sum_{p \in P} (r_{b,p} - \bar{r}_b)^2}}$$

• Possible similarity values between -1 and 1

	Item1	Item2	Item3	Item4	Item5
Alice	5	3	4	4	?
User1	3	1	2	3	3
User2	4	3	4	3	5
User3	3	3	1	5	4
User4	1	5	5	2	1

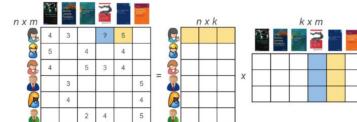
17

💡:

- **Matrix Factorization** is one of the most popular methods for collaborative filtering

- Given rating matrix Y
- Each row represents an user u
- While each column an item i

Regularization of p and q should be applied



$$\hat{y}_{ui} = f(u, i | p_u, q_i) = p_u^T q_i = \sum_{k=1}^K p_{uk} q_{ik}$$

$$L_{sqr} = \sum_{(u,i) \in \mathcal{Y} \cup \mathcal{Y}^-} w_{ui} (y_{ui} - \hat{y}_{ui})^2$$

25

$$pred(a, p) = \bar{r}_a + \frac{\sum_{b \in N} sim(a, b) * (r_{b,p} - \bar{r}_b)}{\sum_{b \in N} |sim(a, b)|}$$

Metrics:

$$MAE = \frac{1}{n} \sum_{i=1}^n |p_i - r_i|$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (p_i - r_i)^2}$$

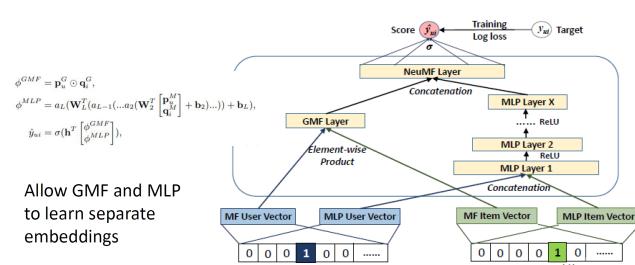
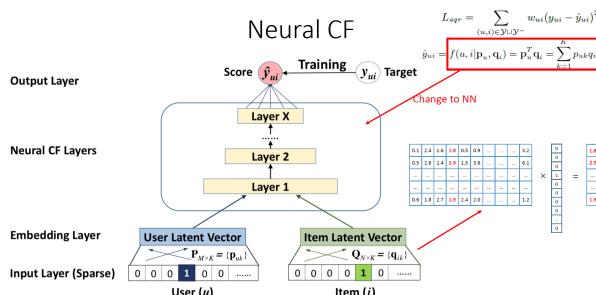
Pros:

- well-understood, works well in some domains, no knowledge engineering required

Cons:

- requires user community, sparsity problems, no integration of other knowledge sources, no explanation of results

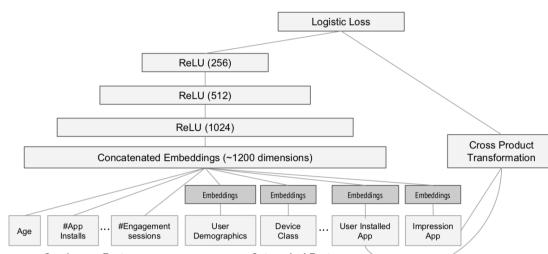
2. Neural Collaborative Filtering (NCF)



32

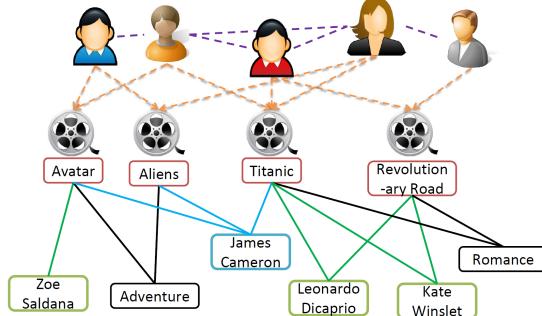
3. Wide & Deep learning

The wide & deep model

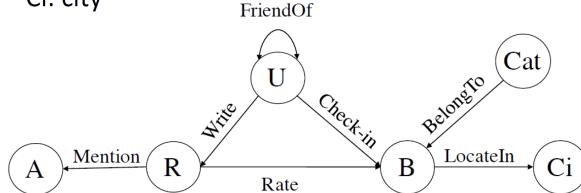


$$P(Y = 1 | \mathbf{x}) = \sigma(\mathbf{w}_{wide}^T \mathbf{x}, \phi(\mathbf{x})] + \mathbf{w}_{deep}^T a^{(lf)} + b)$$

4. Heterogeneous Info for Recommendation System.

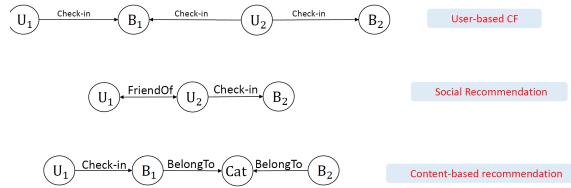


- R: reviews;
- U: users;
- B: business;
- Cat: category of item;
- Ci: city

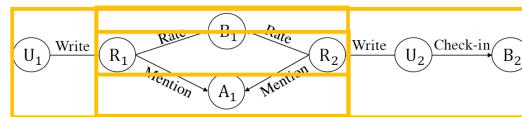


Metapath based RS

- Metapath \rightarrow Recommending Strategy.



Compute Similarity based on a Meta-graph



$$\begin{aligned} \text{Compute } C_{P_1} : & C_{P_1} = W_{RB} \cdot W_{RB}^T; \\ \text{Compute } C_{P_2} : & C_{P_2} = W_{RA} \cdot W_{RA}^T; \\ \text{Compute } C_{S_r} : & C_{S_r} = C_{P_1} \odot C_{P_2}; \\ \text{Compute } C_{M_9} : & C_{M_9} = W_{UR} \cdot C_{S_r} \cdot W_{UR}^T \cdot W_{UB}; \end{aligned}$$

L12 Recommendation Sys 2.

• Problem: DMR doesn't explicitly use graph structure.

• Only First-order graph structure capture in training objective

• Solution: GNN (explicitly, high order)

1. NGCF: NGCF jointly learn:

1. GNN 2. node shallow embedding

- Iteratively update node embeddings using neighboring embeddings.

$$\begin{aligned} h_v^{(k+1)} &= \text{COMBINE}\left(h_v^{(k)}, \text{AGGR}\left(\left\{h_u^{(k)}\right\}_{u \in N(v)}\right)\right) \\ h_u^{(k+1)} &= \text{COMBINE}\left(h_u^{(k)}, \text{AGGR}\left(\left\{h_v^{(k)}\right\}_{v \in N(u)}\right)\right) \end{aligned}$$

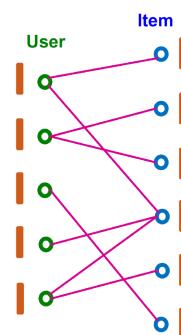
\downarrow
 $\sigma(W \cdot IX, y)$

Mean

• Finally: At k^{th} iteration:

- For all $u \in \mathcal{U}, v \in \mathcal{V}$, we set
 $u \leftarrow h_v^{(K)}, v \leftarrow h_v^{(K)}$.
- Score function is the inner product:
 $\text{Score}(u, v) = u^T v$

• User-item-user: very informative



2. LightGCN.

- Shallow embedding: $O(ND)$
- GNN: $O(D^2)$ $f: R^D \rightarrow R^D$

• Simplifying GCN

$$\cdot \text{GCN: } E^{(k+1)} = \sigma(\tilde{A} E^{(k)} W^{(k)}) \quad , \quad \tilde{A} = D^{-1/2} A D^{-1/2}$$

• Simplify GCN by removing ReLU:

$$\begin{aligned} E^k &= \tilde{A} E^{k-1} W^{k-1} \\ &= \tilde{A} (\tilde{A} E^{k-2}) W^{k-2} W^{k-1} \\ &\dots \\ &= \tilde{A}^k E^0 (W^0 \cdots W^{k-1}) \end{aligned}$$

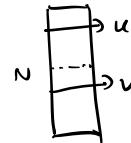
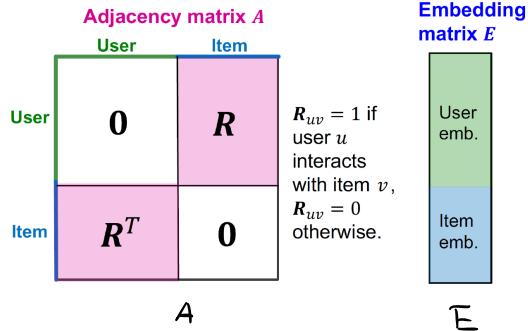
$\underbrace{\qquad\qquad}_{W}$

$$\Rightarrow E^k = \tilde{A}^k E W$$

• Alg: $E \leftarrow \tilde{A} E$ for every k .

$$\begin{aligned} \cdot \text{Final } E: E_{\text{final}} &= \alpha_0 E^0 + \alpha_1 E^1 + \dots + \alpha_k E^k \quad (\alpha = \frac{1}{k+1}) \\ &= \frac{1}{k+1} \sum_{k=0}^k \alpha_k E^k \end{aligned}$$

• E_{final} for score: $u^T v \quad (u, v \in E)$



• Remember: Normalized Laplace: $D^{-1}(I-W)$, $D^{-1}(I-W)D^{-1}$

$$\sum_{k=0}^{\infty} \tilde{A}^k = (I - \tilde{A})^{-1} = \left(D^{-1}(I-W)D^{-1} \right)^{-1} = L^{-1}$$

LightGCN

$$L = U \Lambda U^T, \quad L^{-1} = U \Lambda^{-1} U^T \quad (Lu = \lambda u \Rightarrow L^{-1}Lu = \lambda L^{-1}u \Rightarrow L^{-1}u = \frac{1}{\lambda} u)$$

Low-pass filter

• LightGCN can be thought as a low-pass filter. $(\sum_k \tilde{A}^k) \cdot E = L^{-1} E$

• Simplification leads to better performance. (than NGCF)

3. PinSAGE:

- Learn embeddings z_i : $d(z_{\text{cake1}}, z_{\text{cake2}}) < d(z_{\text{cake1}}, z_{\text{sweater}})$
- Share Negative Sample: V_{neg} for a batch U_{mini} (performance stay similar!)

- Curriculum learning : make negative gradually harder
 - At n-th epoch , Add n-1 hard negative items
- Hard Negatives : close but not connected!
 - method: Run PageRank , random sample from high ranking item (but not too high)