

Machine Learning on Graph

Network Embedding

Presented By

Sanasam Ranbir Singh



Open Source Intelligence Lab
Department of Computer Science and Engineering,
Indian Institute of Technology, Guwahati
Assam, INDIA

<http://www.iitg.ernet.in/cseweb/osint/default/index>

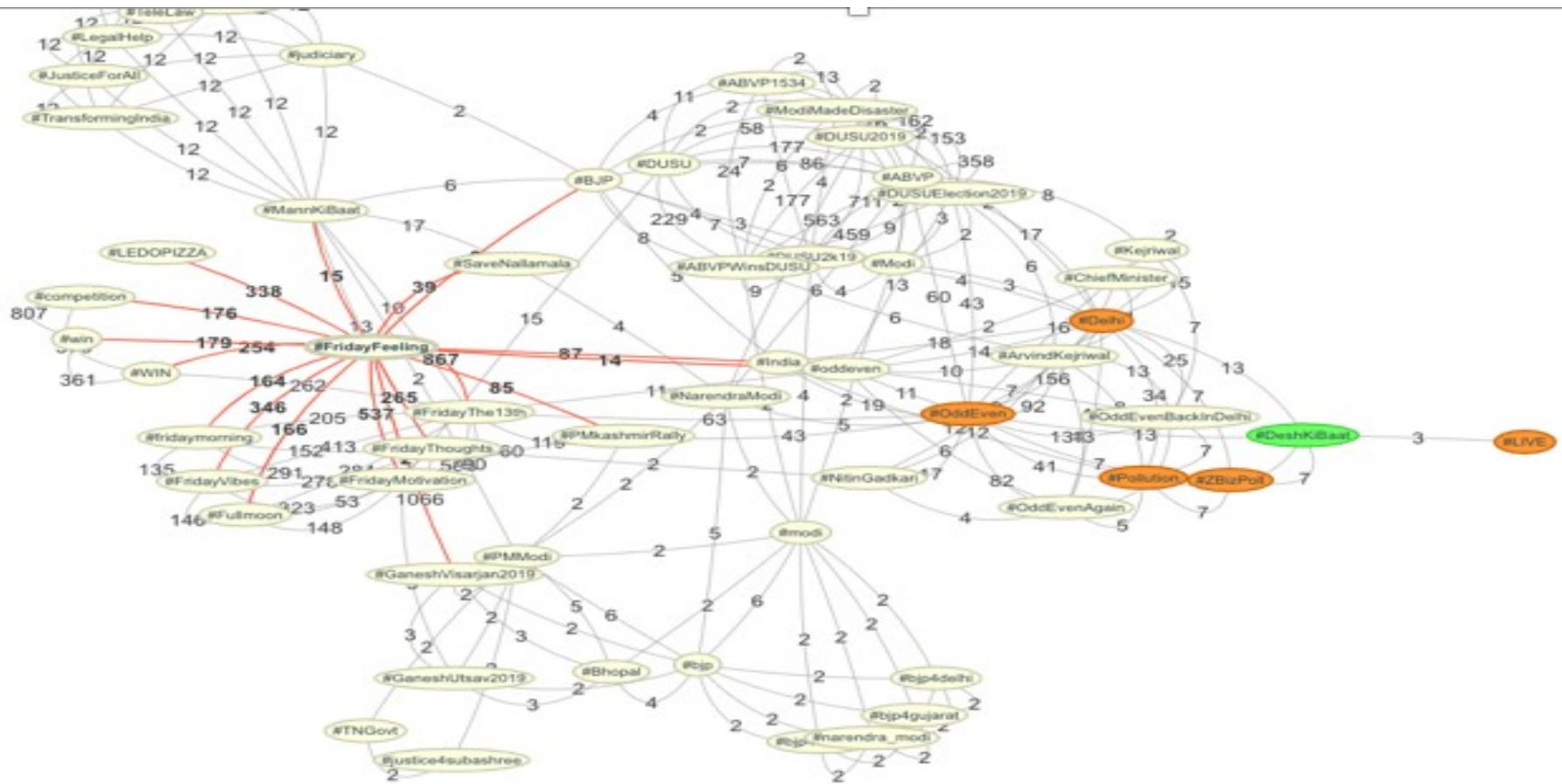
Facebook Friendship Network



Twitter Followers' Network

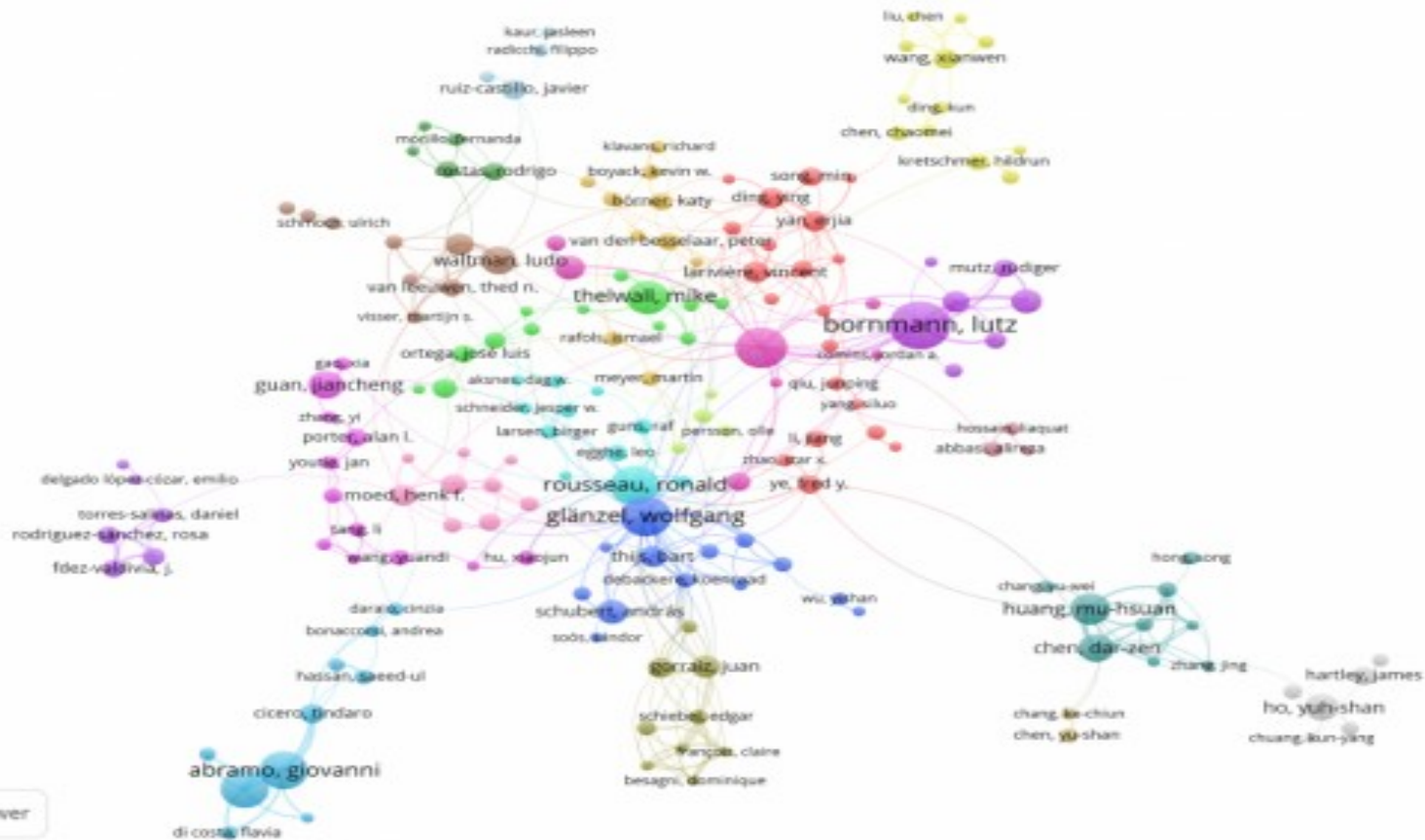


Hashtag Co-occurrence Network



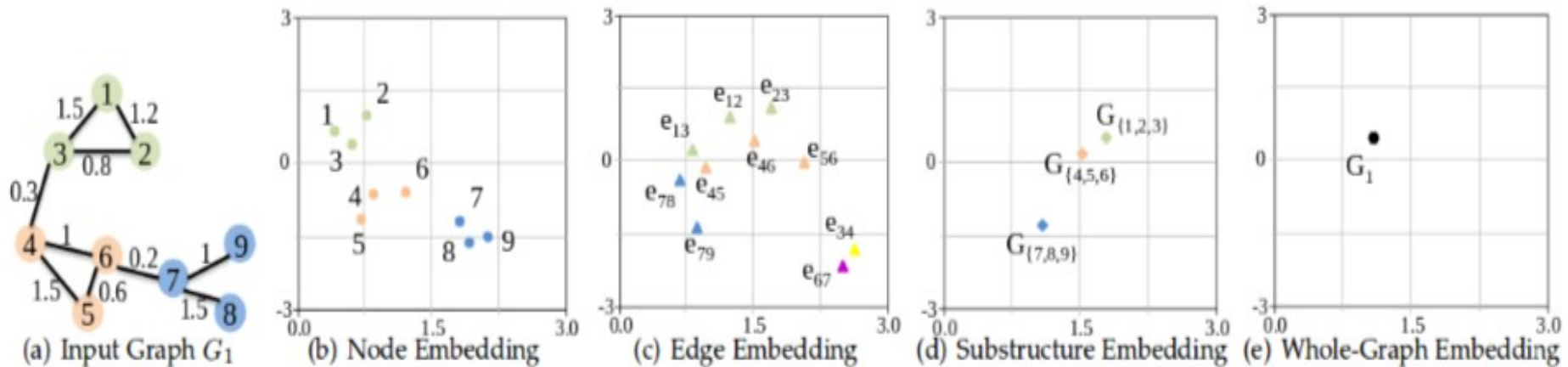
[illegible]

Co-Authorship Network

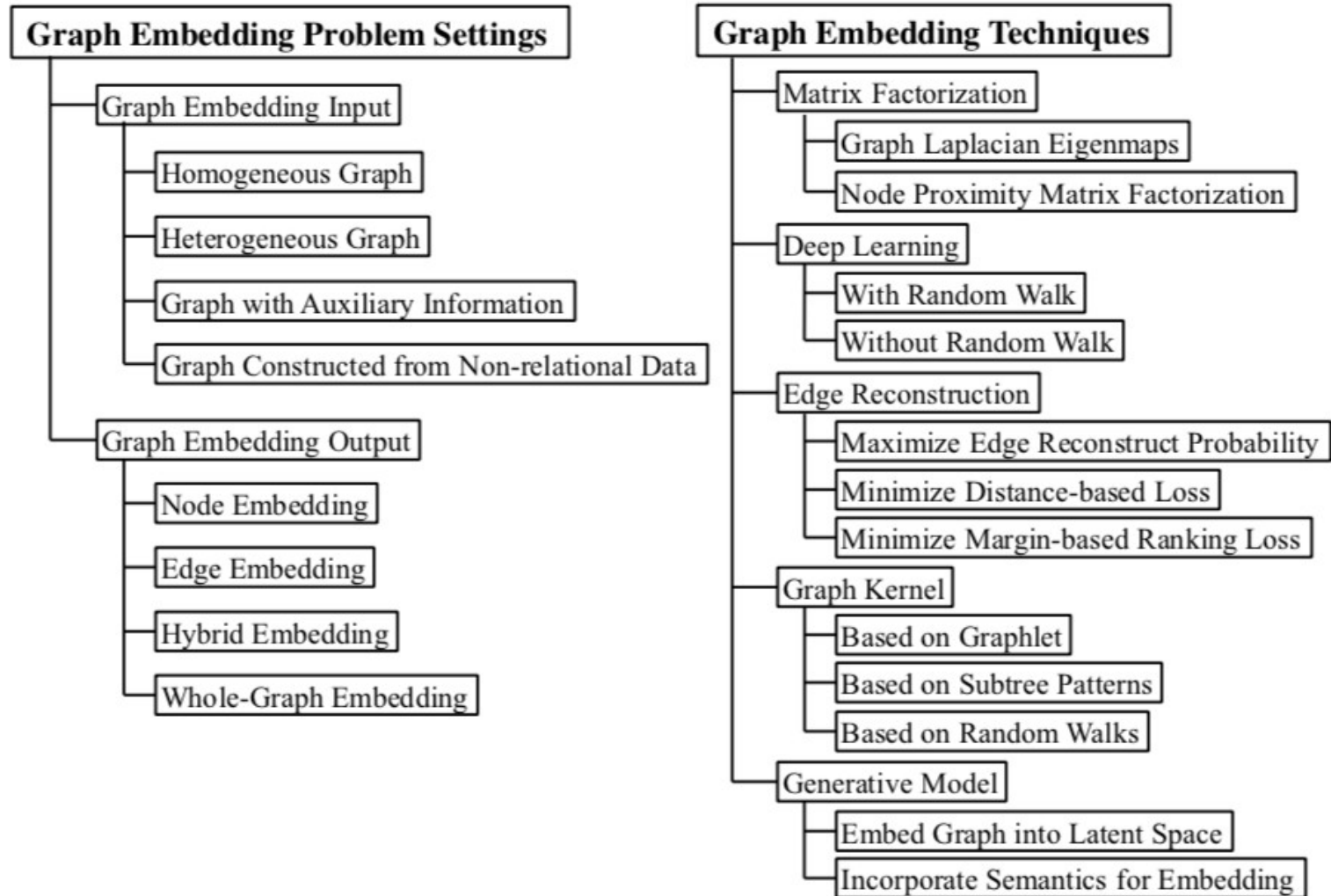


What is Network Embedding?

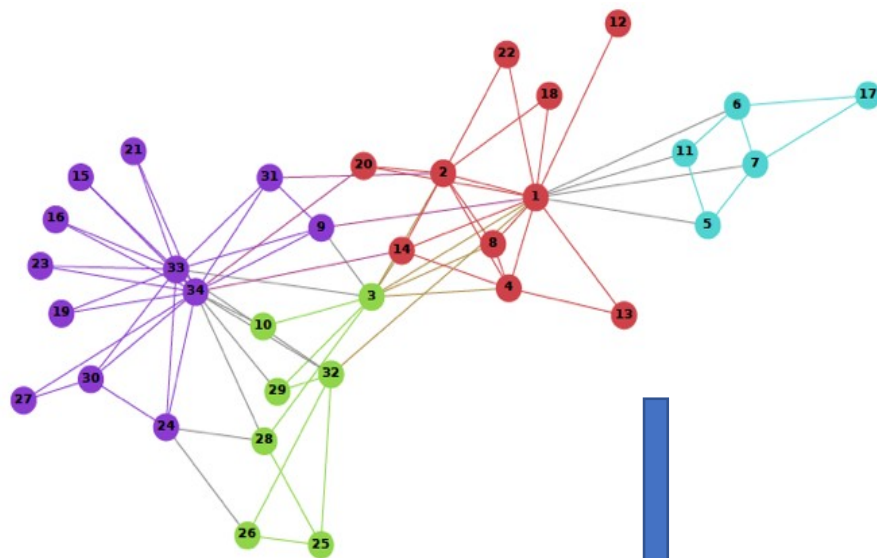
Suppose $G(V,E)$ represents a network then Network Embedding refers to generating low dimensional network features corresponding to Nodes, Edges, Substructures, and the Whole-Graph.



Network Embedding - Texonomy



Node Embedding

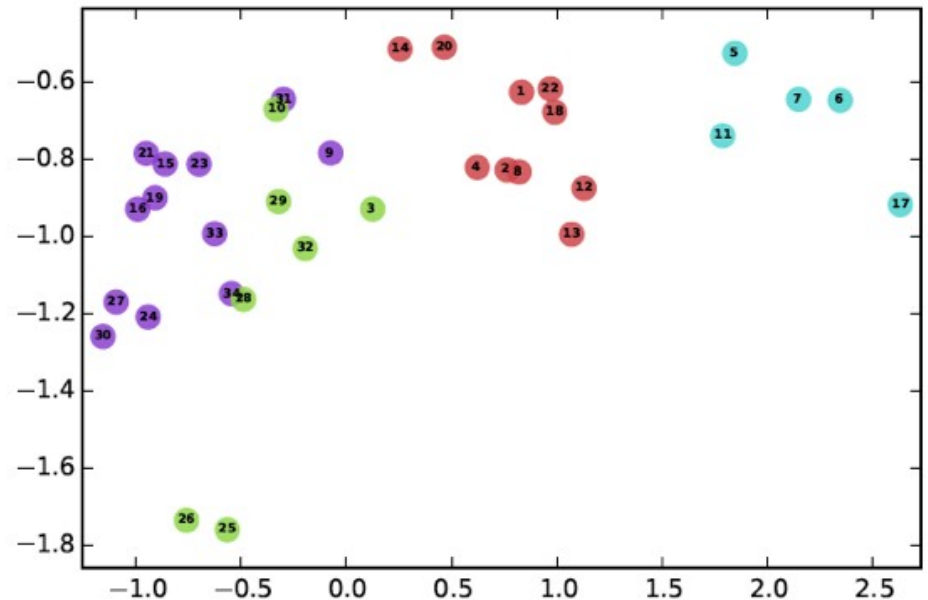
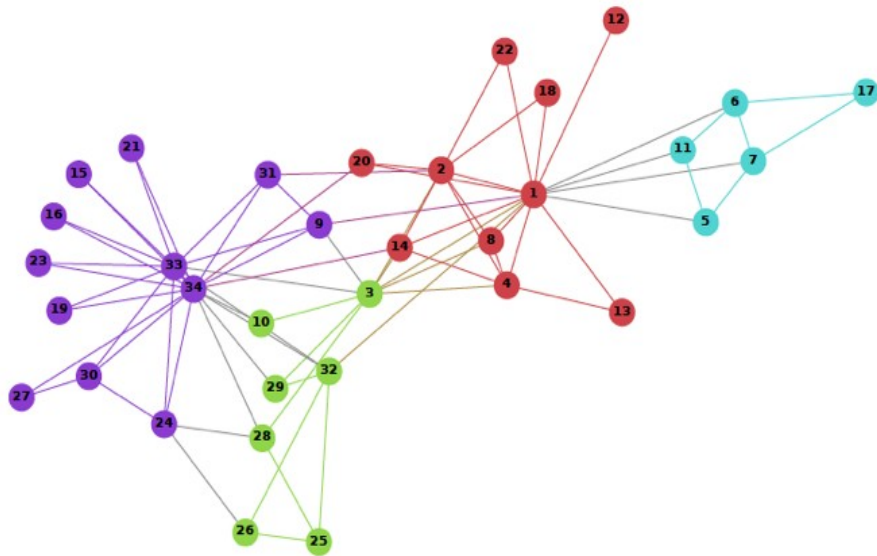


Embedding

Find embedding of nodes to d-dimensions so that “similar” nodes in the graph have embeddings that are close together.

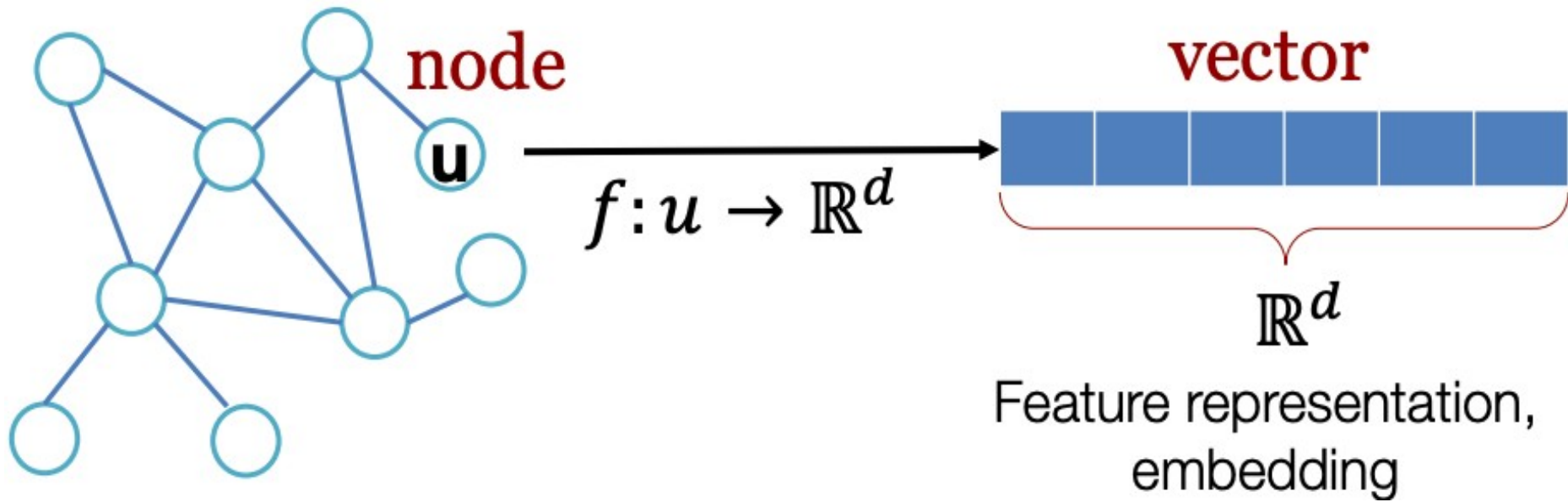
```
0 : [ 1.98479516  2.104902 ]
1 : [ 1.37277632  1.17311989]
2 : [ 0.6859528  1.08576755]
3 : [ 1.27303975  0.41918123]
4 : [ 0.65722556 -0.36297315]
5 : [ 0.71494817 -0.39592245]
6 : [ 0.71494817 -0.39592245]
7 : [ 1.04116841 -0.05324659]
8 : [-0.27647113 -0.24436752]
9 : [ 0.65722556 -0.36297315]
10 : [ 0.37478198 -0.46299671]
11 : [ 0.63250419 -0.30641371]
12 : [ 0.66706122 -0.24109217]
13 : [ 0.64664541 -0.38866438]
14 : [ 0.27253823 -0.57650996]
15 : [ 0.64664541 -0.38866438]
16 : [-0.5209427 -0.58291863]
17 : [-0.48328386 -0.17530251]
18 : [-0.24844813 -0.38526576]
19 : [-0.5175929 -0.23180349]
20 : [-0.34294104 -0.09404636]
21 : [-1.56457792  1.8504155 ]
22 : [ 0.28271458 -0.28822466]
23 : [-1.76282103  3.02650113]
24 : [-0.70912996 -0.57803913]
25 : [-0.70912996 -0.57803913]
26 : [-0.70912996 -0.57803913]
27 : [-0.70912996 -0.57803913]
28 : [-0.70912996 -0.57803913]
29 : [-1.08578976 -0.27525378]
30 : [-0.37477941  0.06842669]
31 : [-1.03412988 -0.27920574]
32 : [-0.28044074  0.04244915]
33 : [-0.58710262 -0.38879992]
```

Node Embedding



Find embedding of nodes to d-dimensions so that “similar” nodes in the graph have embeddings that are close together.

Node Embedding



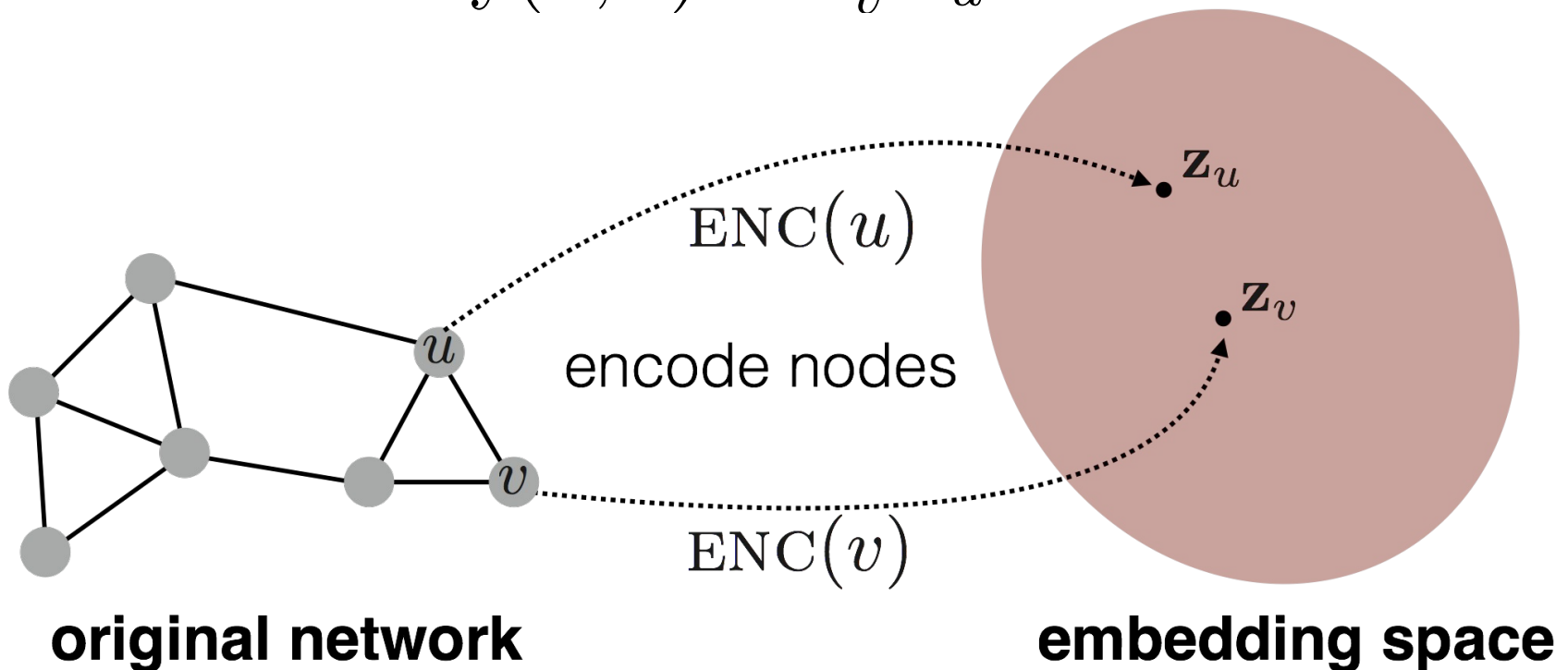
Find embedding of nodes to d-dimensions so that “similar” nodes in the graph have embeddings that are close together.

Node Embedding

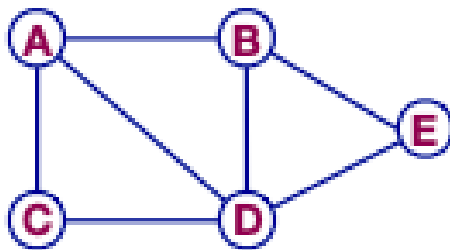
Goal is to encode nodes so that similarity in the embedding space (e.g., dot product)

approximates similarity in the original network.

$$\text{similarity}(u, v) \approx \mathbf{z}_v^\top \mathbf{z}_u$$



Vectorization of Nodes



	A	B	C	D	E	Feat	
A	0	1	1	1	0	1	0
B	1	0	0	1	1	0	0
C	1	0	0	1	0	0	1
D	1	1	1	0	1	1	1
E	0	1	0	1	0	1	0

Many possible ways to create node features:

- Node degree, PageRank score, motifs, ...
- Degree of neighbors, PageRank of neighbors, ...

Learning Node Embedding

1. Define an encoder :

$$\text{ENC}(v) = \mathbf{z}_v$$

d-dimensional embedding

2. Define a node similarity function in the original network:
 $\text{similarity}(u, v)$

Similarity of u and v in the original network

3. Optimize the parameters of the encoder so that:

$$\text{similarity}(u, v) \approx \mathbf{z}_v^\top \mathbf{z}_u$$

Encoding

Once we find encoding matrix \mathbf{Z} , encoding of a node v can be defined as

$$\text{ENC}(v) = \mathbf{Z}\mathbf{v}$$

$\mathbf{Z} \in \mathbb{R}^{d \times |\mathcal{V}|}$ matrix, each column is node embedding

$\mathbf{v} \in \mathbb{I}^{|\mathcal{V}|}$ indicator vector, all zeroes except a one in column indicating node v

Encoding

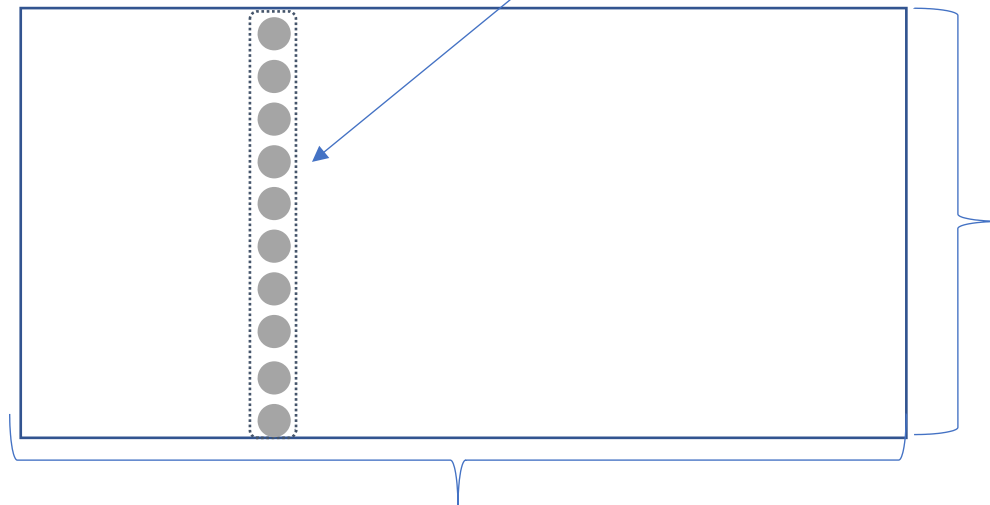
Once we find encoding matrix \mathbf{Z} , encoding of a node v (one hot vector) can be defined as

$$\text{ENC}(v) = \mathbf{Z}\mathbf{v}$$

embedding vector for
a specific node

embedding
matrix

$\mathbf{Z} =$



Dimension/
size of
embeddings

one column per

Learning Node Embedding

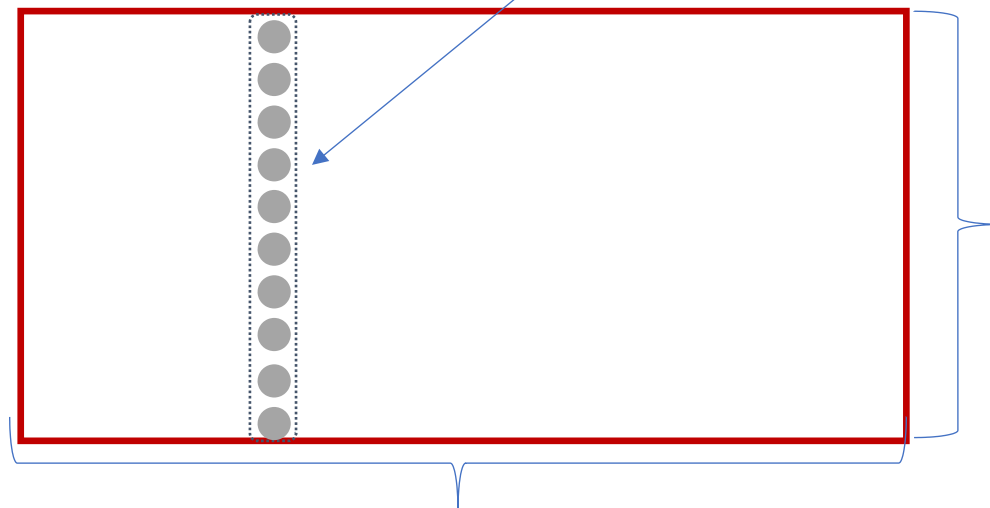
Learning node embedding is to learn the matrix \mathbf{Z}

$$\text{ENC}(v) = \mathbf{Z}\mathbf{v}$$

embedding vector for
a specific node

embedding
matrix

$$\mathbf{Z} =$$



Dimension/
size of
embeddings

one column per

Three ways of Learning Z

1. Adjacency-based similarity
2. Multi-hop similarity
3. Deep Learning
 - Unsupervised Approach(Random walk)
 - DeepWalk
 - Supervised Approach
 - GNN
 - GraphSAGE
 - Graph-BERT
 - GCN

Adjacency-based similarity

Let us assume that A is the adjacency matrix and two nodes are connected if they are same

Then, the loss function is

$$\mathcal{L} = \sum_{(u,v) \in V \times V} \| \mathbf{z}_u^\top \mathbf{z}_v - A_{u,v} \|^2$$

loss (what we want to minimize)

sum over all node pairs

embedding similarity

(weighted) adjacency matrix for the graph

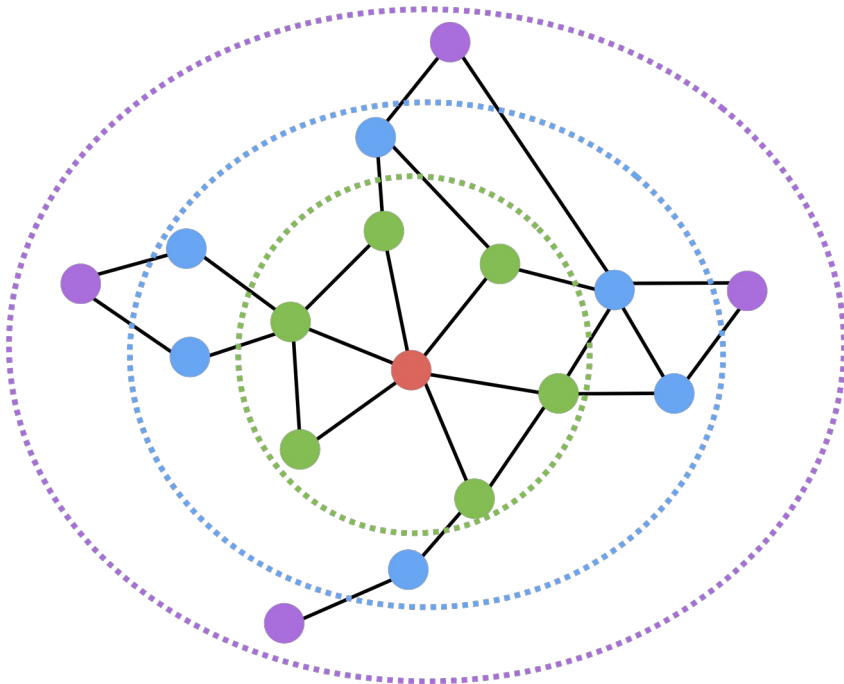
Two approaches

$$\mathcal{L} = \sum_{(u,v) \in V \times V} \|\mathbf{z}_u^\top \mathbf{z}_v - \mathbf{A}_{u,v}\|^2$$

- Find embedding matrix that minimizes the loss
 - Option 1: Stochastic gradient descent (SGD).
 - Option 2: Matrix decomposition(SVD).

Multi-Hop Similarity

- **Idea:** Consider k-hop node neighbors.
 - E.g., two or three-hop neighbors.



- **Red:** Target node
- **Green:** 1-hop neighbors
 - A (i.e., adjacency matrix)
- **Blue:** 2-hop neighbors
 - A^2
- **Purple:** 3-hop neighbors
 - A^3

Katz link prediction method

Multi-Hops Similarity Embedding

$$\mathcal{L} = \sum_{(u,v) \in V \times V} \left\| \mathbf{z}_u^\top \mathbf{z}_v - \mathbf{S}_{u,v} \right\|^2$$

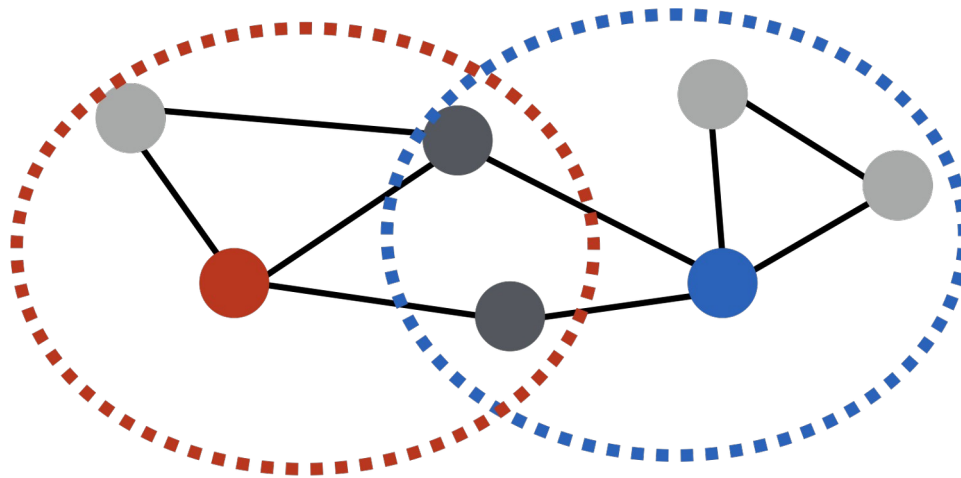
embedding similarity

multi-hop network similarity (i.e., any neighborhood overlap measure)

$\mathbf{S}_{u,v}$ is the neighborhood overlap between u and v .

HOPE (Yan et al., 2016)

- **Another option:** Measure overlap between node neighborhoods.



- Example overlap functions:
 - Common Neighbours
 - Jaccard Co-efficient
 - Adamic-Adar score

GraRep from Cao et al, 2015

- **Basic idea:**
$$\mathcal{L} = \sum_{(u,v) \in V \times V} \|\mathbf{z}_u^\top \mathbf{z}_v - \mathbf{A}_{u,v}^k\|^2$$

- Train embeddings to predict k-hop neighbors.

- Probabilistic Adjacency Matrix:

- Use log-transformed, probabilistic adjacency matrix:

$$\tilde{\mathbf{A}}_{i,j}^k = \max \left(\log \left(\frac{(\mathbf{A}_{i,j}/d_i)}{\sum_{l \in V} (\mathbf{A}_{l,j}/d_l)^k} \right)^k - \alpha, 0 \right)$$

node degree constant

- Train multiple different hop lengths and concatenate output.

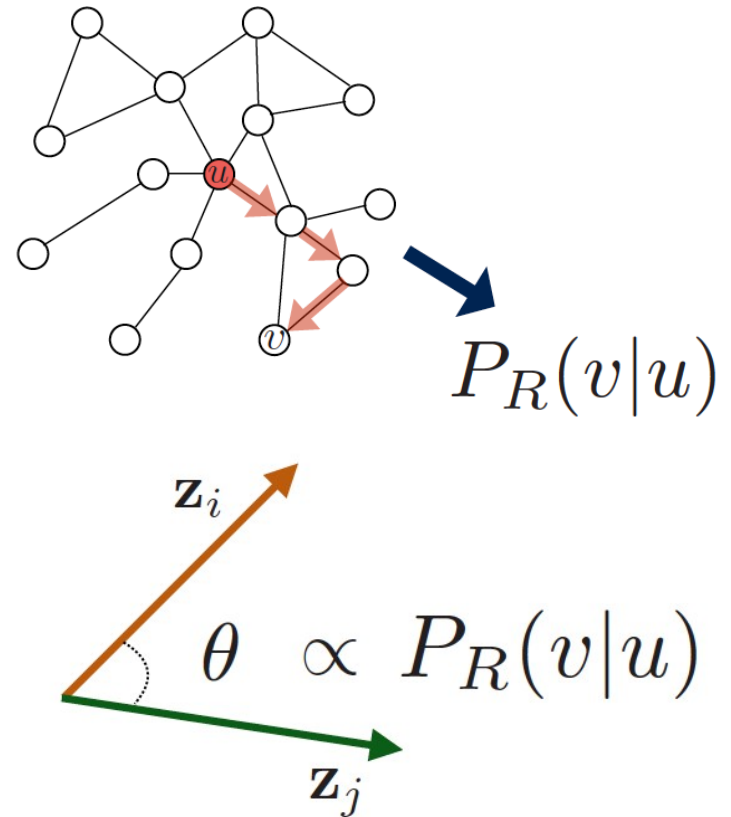
Unsupervised: Random Walk Approach

$$\mathbf{z}_u^\top \mathbf{z}_v \approx$$

probability that u
and v co-occur
on a random
walk over the
network

Random Walk Approach

1. Estimate probability of visiting node v on a random walk starting from node u using some random walk strategy R .
2. Optimize embeddings to encode these random walk statistics.



Random Walk Optimization

1. Run short random walks starting from each node on the graph using some strategy R .
2. For each node u collect $N_R(u)$ node sequence visited on random walks starting from u .
3. Optimize embeddings to according to:

$$\mathcal{L} = \sum_{u \in V} \sum_{v \in N_R(u)} -\log(P(v|\mathbf{z}_u))$$

Random Walk Optimization

$$\mathcal{L} = \sum_{u \in V} \sum_{v \in N_R(u)} -\log(P(v|\mathbf{z}_u))$$

- **Intuition:** Optimize embeddings to maximize likelihood of random walk co-occurrences.
- **Parameterize** $P(v|\mathbf{z}_u)$ **using softmax:**
$$P(v|\mathbf{z}_u) = \frac{\exp(\mathbf{z}_u^\top \mathbf{z}_v)}{\sum_{n \in V} \exp(\mathbf{z}_u^\top \mathbf{z}_n)}$$

Random Walk Optimization

Putting things together:

$$\mathcal{L} = \sum_{u \in V} \sum_{v \in N_R(u)} -\log \left(\frac{\exp(\mathbf{z}_u^\top \mathbf{z}_v)}{\sum_{n \in V} \exp(\mathbf{z}_u^\top \mathbf{z}_n)} \right)$$

sum over all nodes u sum over nodes v seen on random walks starting from u predicted probability of u and v co-occurring on random walk

Optimizing random walk embeddings =

Finding embeddings \mathbf{z} that

Random Walk Optimization

Solution: Negative sampling

$$\log \left(\frac{\exp(\mathbf{z}_u^\top \mathbf{z}_v)}{\sum_{n \in V} \exp(\mathbf{z}_u^\top \mathbf{z}_n)} \right)$$

$$\approx \log(\sigma(\mathbf{z}_u^\top \mathbf{z}_v)) - \sum_{i=1}^k \log(\sigma(\mathbf{z}_u^\top \mathbf{z}_{n_i})), n_i \sim P_V$$

sigmoid
function

random
distribution over
all nodes

i.e., instead of normalizing w.r.t. all nodes,
just normalize against k random
“negative samples”

Random Walk

- Just run fixed-length, unbiased random walks starting from each node (i.e., [DeepWalk, Perozzi et al., 2013](#)).
- Use flexible, biased random walks that can trade off between **local** and **global** views of the network (i.e., [Node2Vec, Grover and Leskovec, 2016](#))

How should we randomly walk?

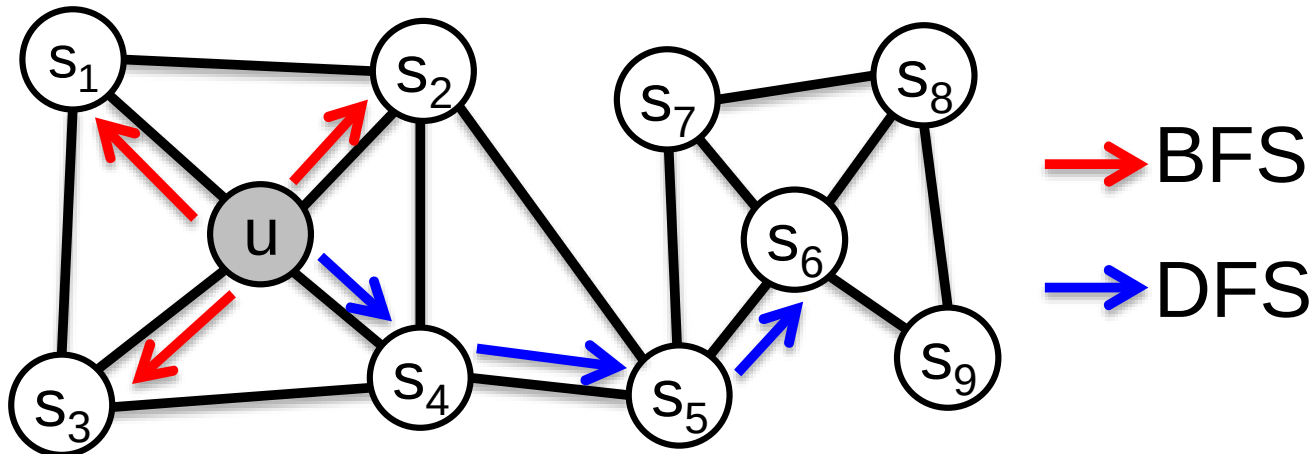
- So far we have described how to optimize embeddings given random walk statistics.
- **What strategies should we use to run these random walks?**
 - Simplest idea: **Just run fixed-length, unbiased random walks starting from each node** (i.e., [DeepWalk from Perozzi et al., 2013](#)).
 - But can we do better?

Node2vec

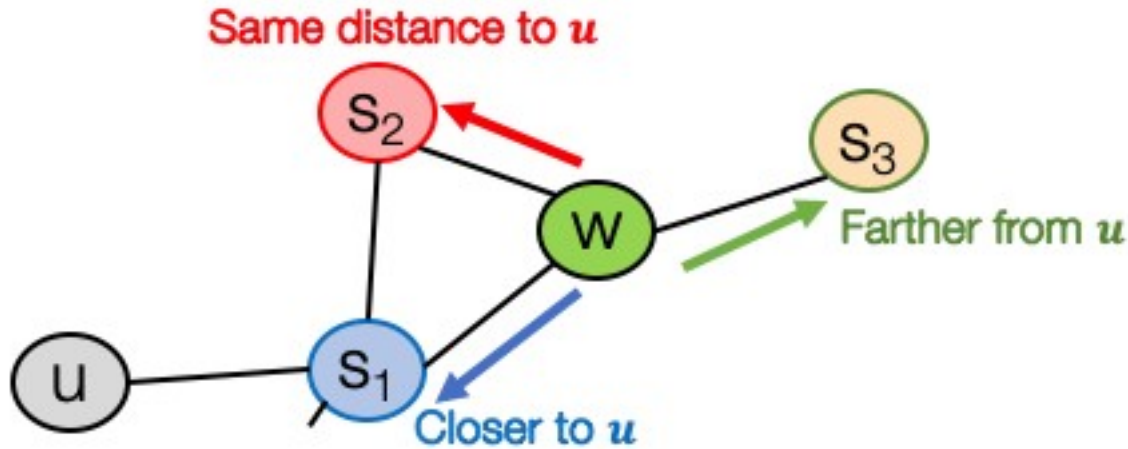
Interpolating BFS and DFS

$$N_{BFS}(u) = \{s_1, s_2, s_3\}$$

$$N_{DFS}(u) = \{s_4, s_5, s_6\}$$



Node2vec: two parameters



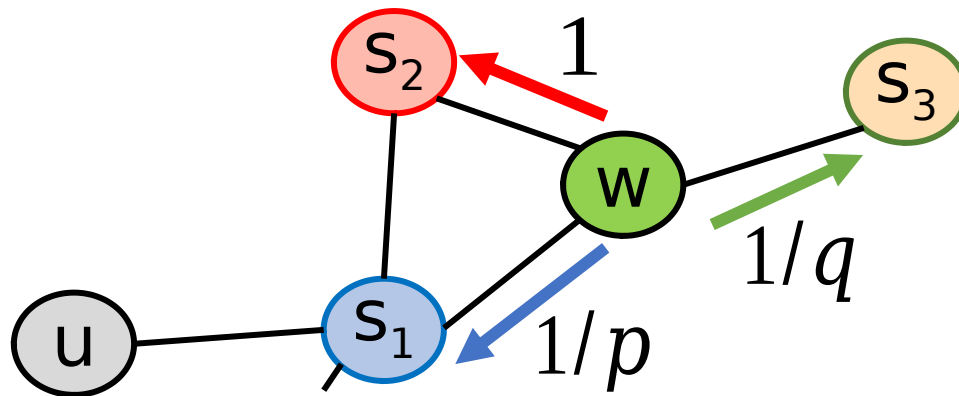
Interpolating BFS and DFS

Biased random walk that given a node generates neighborhood

- Two parameters:
 - Return parameter :
 - Return back to the previous node
 - In-out parameter :
 - Moving outwards (DFS) vs. inwards (BFS)

Node2vec: two parameters

- Walker is at w . Where to go next?

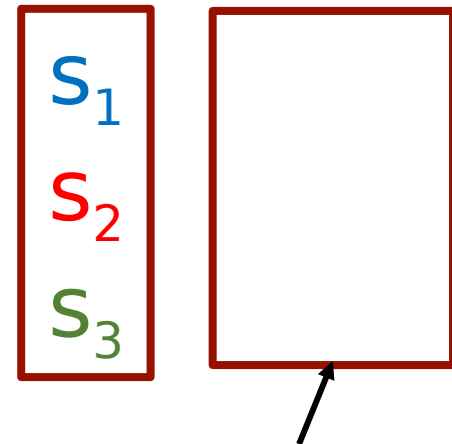
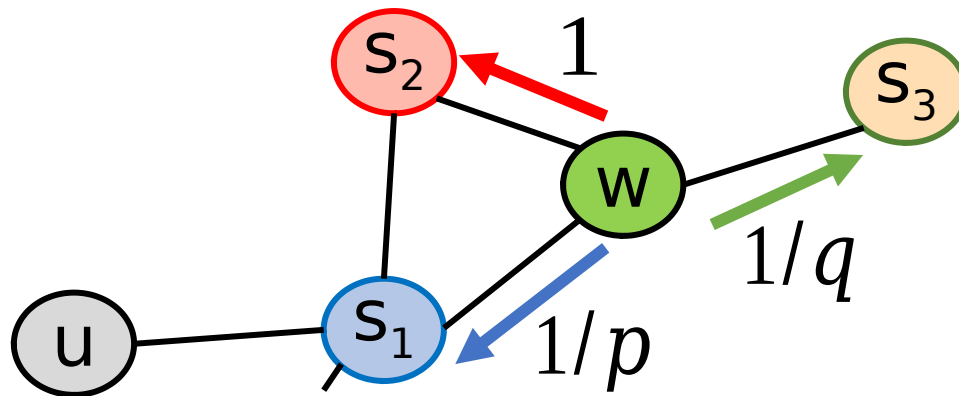


are
unnormaliz
ed
probabilitie
s

- model transition probabilities
 - ... return parameter
 - ... "walk away" parameter

Node2vec: two parameters

- Walker is at w . Where to go next?

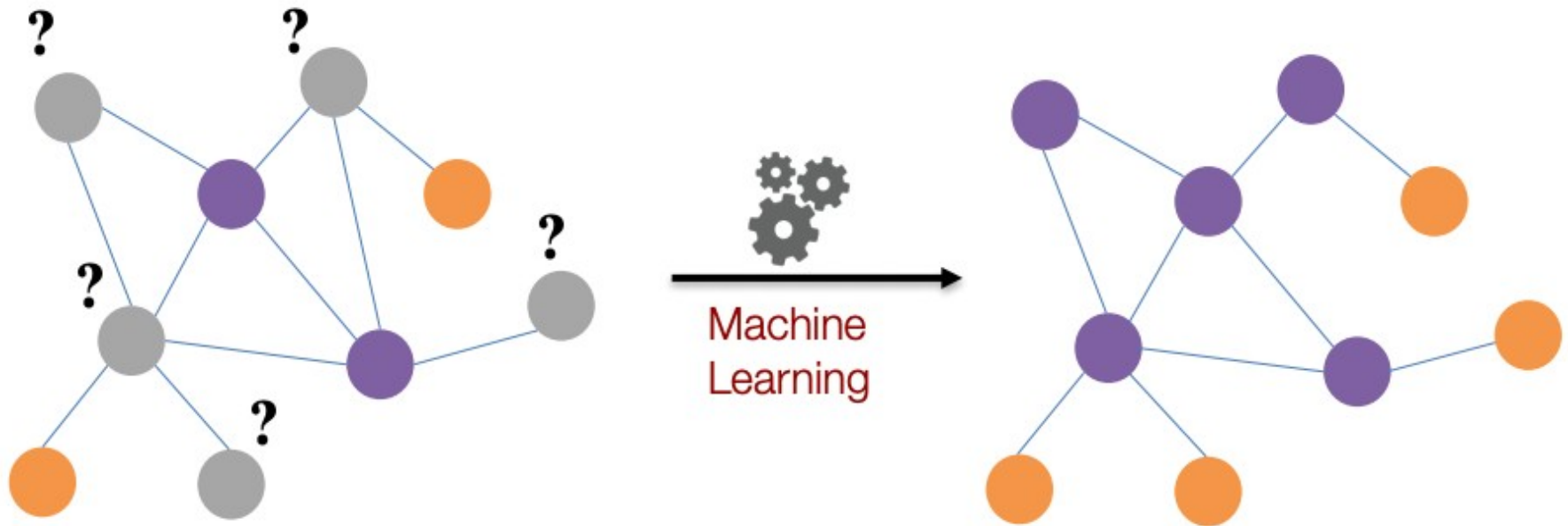


Unnormalized
transition prob.

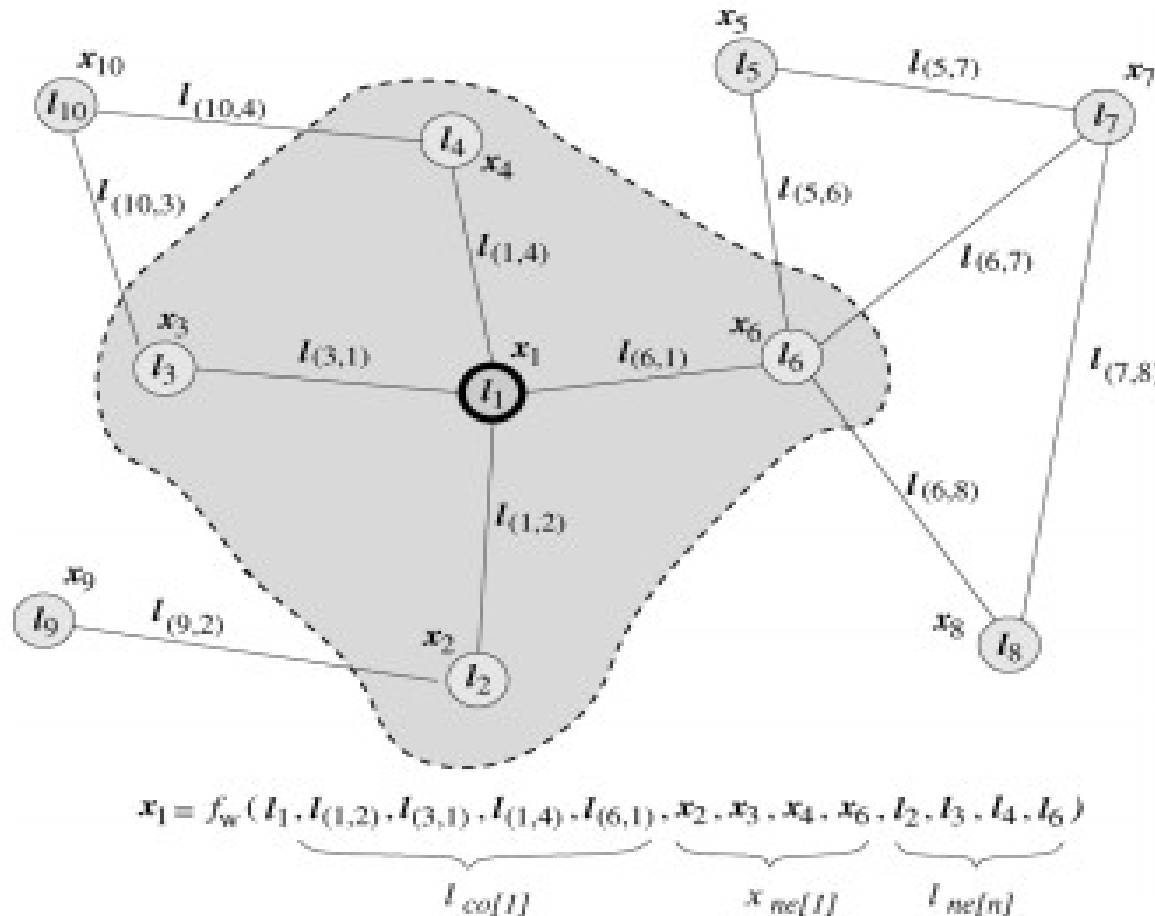
- **BFS-like** walk: Low value of
- **DFS-like** walk: Low value of

are the nodes visited by the walker

Supervised Approach



Graph Neural Network



$$\mathbf{x}_n = f_{\mathbf{w}}(\mathbf{l}_n, \mathbf{l}_{\text{cof}[n]}, \mathbf{x}_{\text{ne}[n]}, \mathbf{l}_{\text{ne}[n]})$$

$$\mathbf{o}_n = g_{\mathbf{w}}(\mathbf{x}_n, \mathbf{l}_n)$$

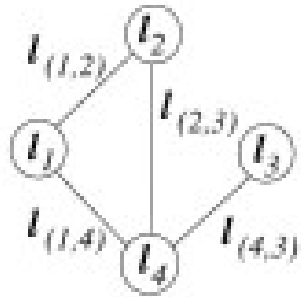
$$\mathbf{x} = F_{\mathbf{w}}(\mathbf{x}, \mathbf{l})$$

$$\mathbf{o} = G_{\mathbf{w}}(\mathbf{x}, \mathbf{l}_N)$$

$$\text{loss} = \sum_{i=1}^P (\mathbf{t}_i - \mathbf{o}_i)$$

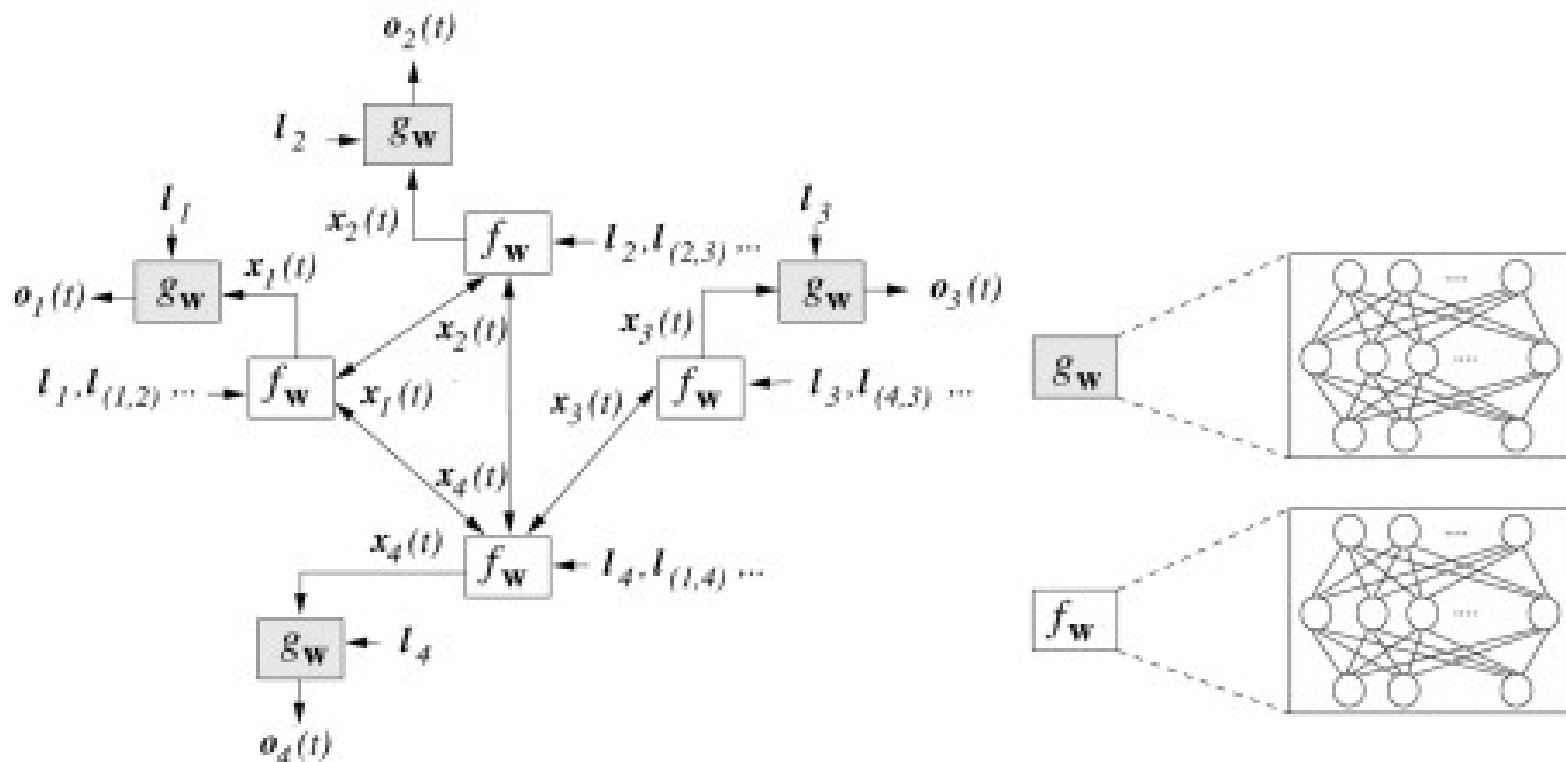
F. Scarselli, M. Gori, A. C. Tsoi, M. Hagenbuchner, and G. Monfardini, "The graph neural network model," IEEE TNN 2009, vol. 20, no. 1, pp. 61–80, 2009.

Graph Neural Network

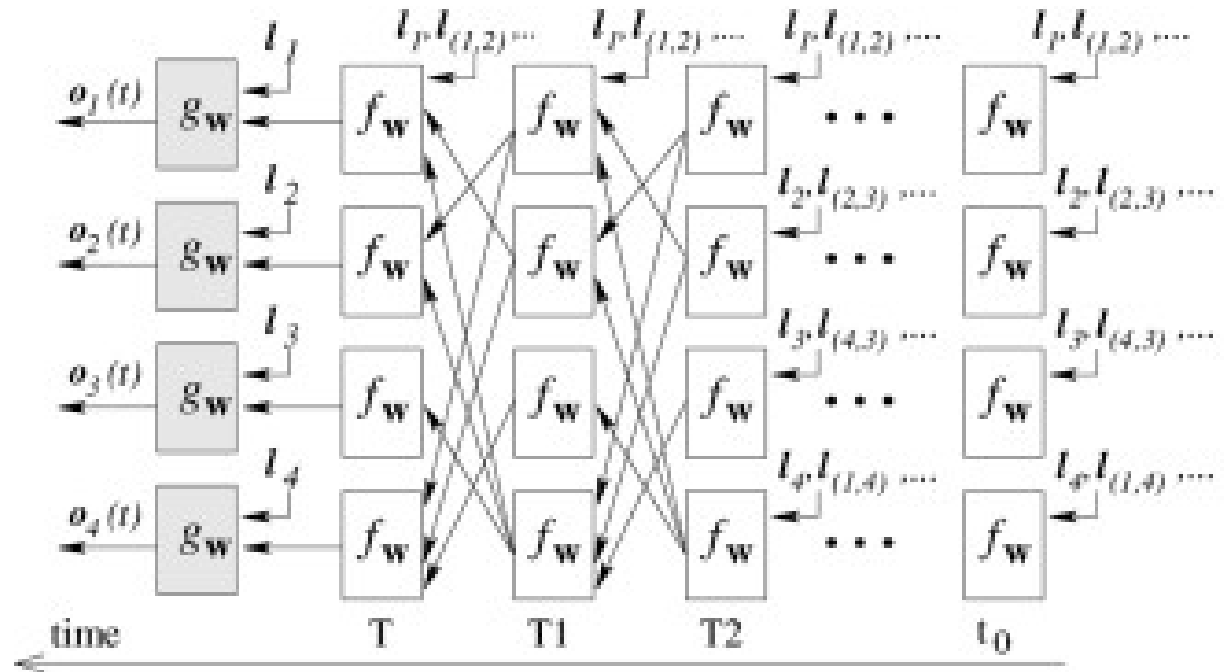
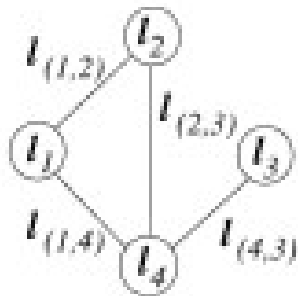


$$\mathbf{x}_n = f_{\mathbf{w}}(\mathbf{l}_n, \mathbf{l}_{\text{co}[n]}, \mathbf{x}_{\text{ne}[n]}, \mathbf{l}_{\text{ne}[n]})$$

$$\mathbf{o}_n = g_{\mathbf{w}}(\mathbf{x}_n, \mathbf{l}_n)$$



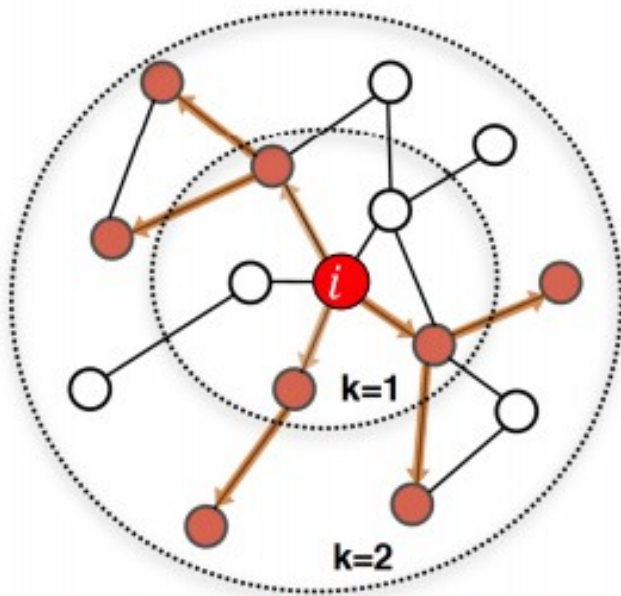
Graph Neural Network



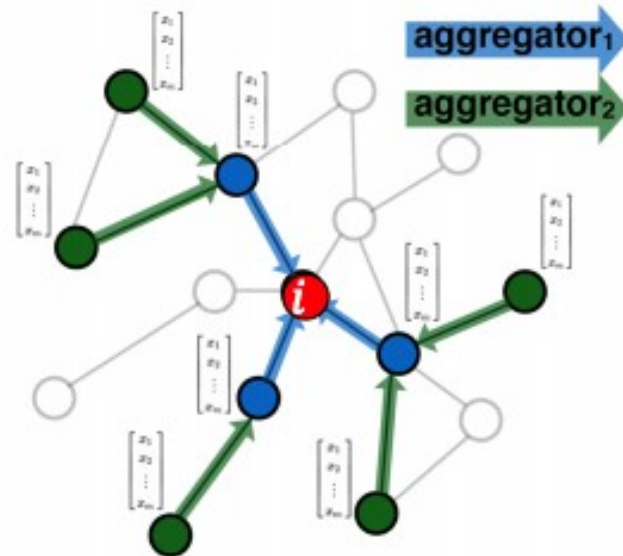
$$loss = \sum_{i=1}^p (t_i - o_i)$$

GraphSAGA

Idea: Node's neighborhood defines a computation graph



Determine node
computation graph



Propagate and
transform information

W. L. Hamilton, Z. Ying, and J. Leskovec, "Inductive representation learning on large graphs," NIPS 2017, pp. 1024–1034, 2017

GraphSAGA

Algorithm 1: GraphSAGE embedding generation (i.e., forward propagation) algorithm

Input : Graph $\mathcal{G}(\mathcal{V}, \mathcal{E})$; input features $\{\mathbf{x}_v, \forall v \in \mathcal{V}\}$; depth K ; weight matrices $\mathbf{W}^k, \forall k \in \{1, \dots, K\}$; non-linearity σ ; differentiable aggregator functions $\text{AGGREGATE}_k, \forall k \in \{1, \dots, K\}$; neighborhood function $\mathcal{N} : v \rightarrow 2^{\mathcal{V}}$

Output : Vector representations \mathbf{z}_v for all $v \in \mathcal{V}$

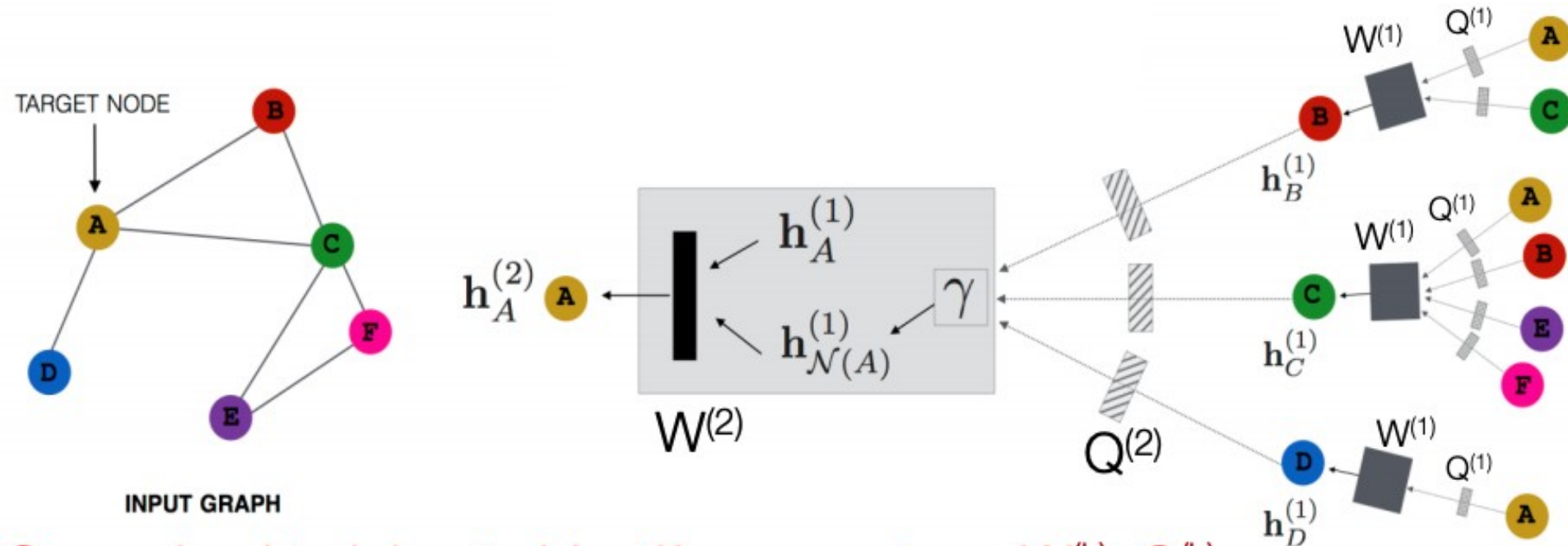
```
1  $\mathbf{h}_v^0 \leftarrow \mathbf{x}_v, \forall v \in \mathcal{V}$  ;  
2 for  $k = 1 \dots K$  do  
3   for  $v \in \mathcal{V}$  do  
4      $\mathbf{h}_{\mathcal{N}(v)}^k \leftarrow \text{AGGREGATE}_k(\{\mathbf{h}_u^{k-1}, \forall u \in \mathcal{N}(v)\})$ ;  
5      $\mathbf{h}_v^k \leftarrow \sigma \left( \mathbf{W}^k \cdot \text{CONCAT}(\mathbf{h}_v^{k-1}, \mathbf{h}_{\mathcal{N}(v)}^k) \right)$   
6   end  
7    $\mathbf{h}_v^k \leftarrow \mathbf{h}_v^k / \|\mathbf{h}_v^k\|_2, \forall v \in \mathcal{V}$   
8 end  
9  $\mathbf{z}_v \leftarrow \mathbf{h}_v^K, \forall v \in \mathcal{V}$ 
```

Update for node i :

$$\mathbf{h}_i^{(k+1)} = \text{ReLU} \left(\mathbf{W}^{(k)} \mathbf{h}_i^{(k)}, \sum_{n \in \mathcal{N}(i)} \left(\text{ReLU}(\mathbf{Q}^{(k)} \mathbf{h}_n^{(k)}) \right) \right)$$

- $\mathbf{h}_i^{(0)} = \mathbf{X}_i$ (directly leverage node attributes)
- $\Sigma(\cdot)$: Aggregator function (avg., LSTM, max-pooling)

GraphSAGA



Supervised training to identify parameters: $W^{(k)}$, $Q^{(k)}$

