

# Graph Embedding

Prof. O-Joun Lee

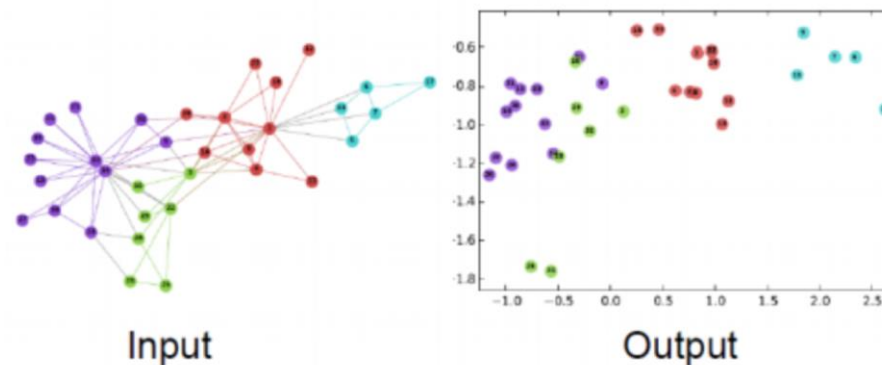
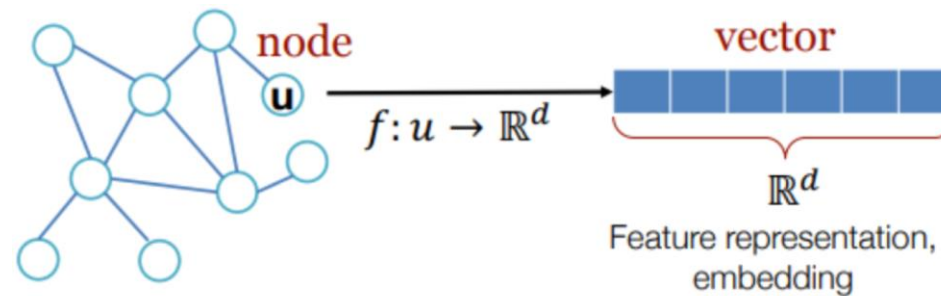
Dept. of Artificial Intelligence,  
The Catholic University of Korea  
*ojlee@catholic.ac.kr*

# Contents



- Graph representation learning
- Representative models
  - Node2Vec
  - LINE
  - SDNE
  - Subgraph2Vec

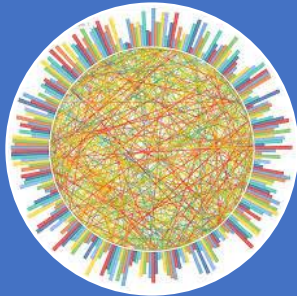
- The objective of a graph embedding: **nodes which are similar** in the graph, should be mapped **close in the vector space**
- Schematic of graph (node) embedding





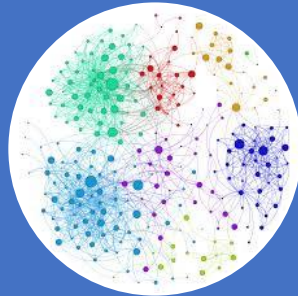
## Network Compression

- A compression store networks more efficiently run graph algorithms faster



## Visualization

- View it from an information visualization perspective



## Community Mining

- Network partitioning ( k-means on the embedding to cluster the nodes )



## Link Prediction

- Predict either missing interactions or links that may appear in the future



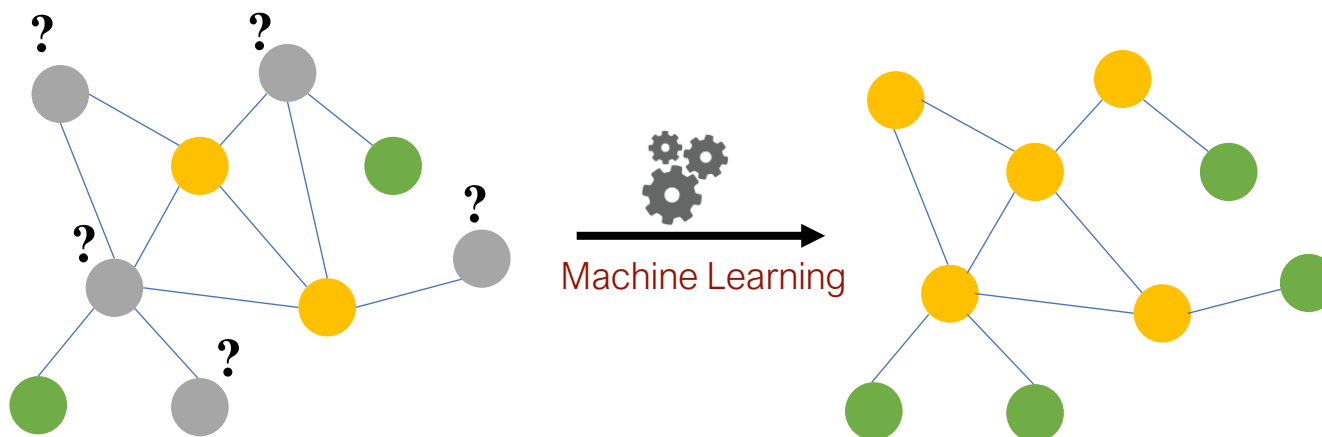
## Node Classification

- Label the full graph based only on this small initial seed set

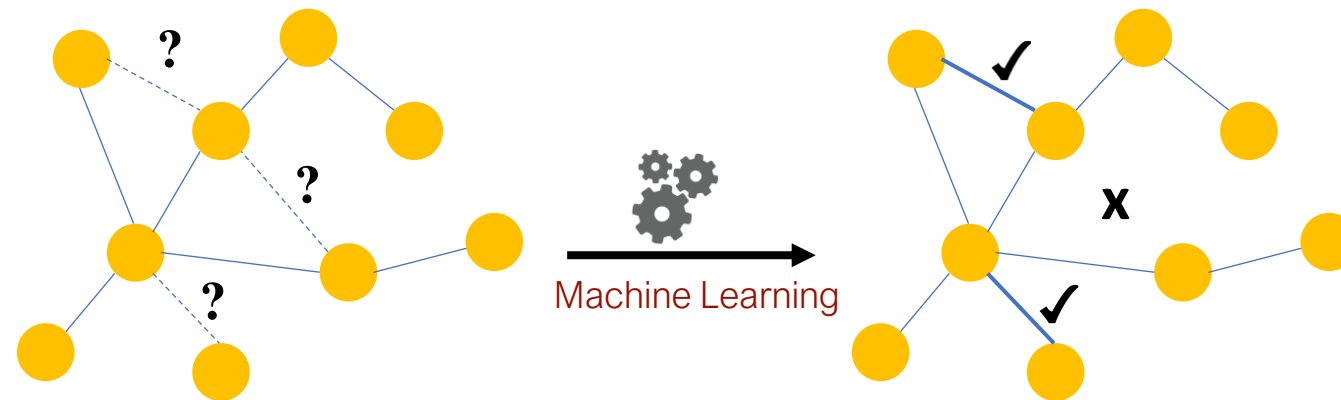




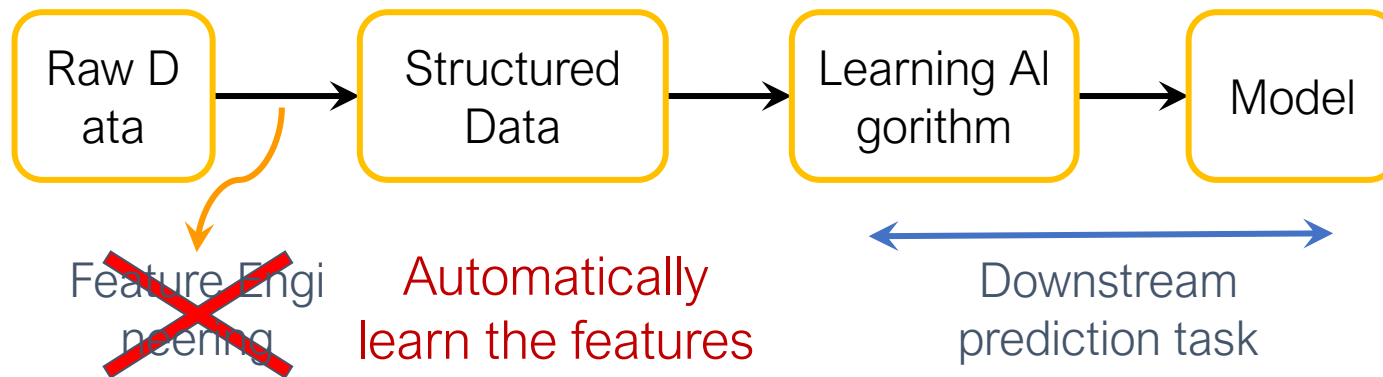
- **Node Classification** is a machine learning task in graph-based data analysis, where the goal is to assign labels to nodes in a graph based on the properties



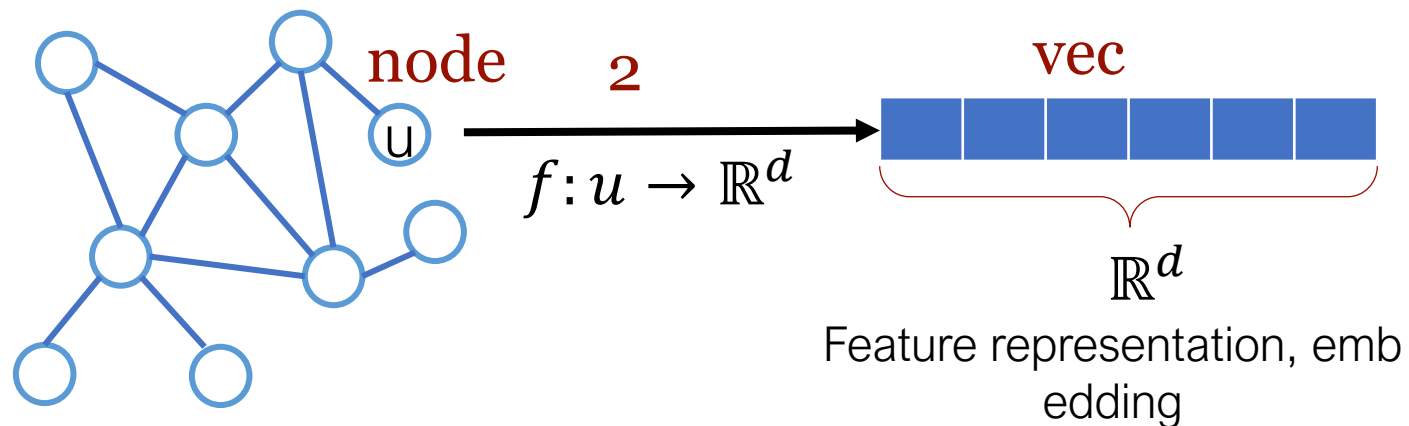
- **Link Prediction** is a task in graph and network analysis where the goal is to predict missing or future connections between nodes in a network.



- (Supervised) Machine Learning Lifecycle: This feature, that feature. Every single time!

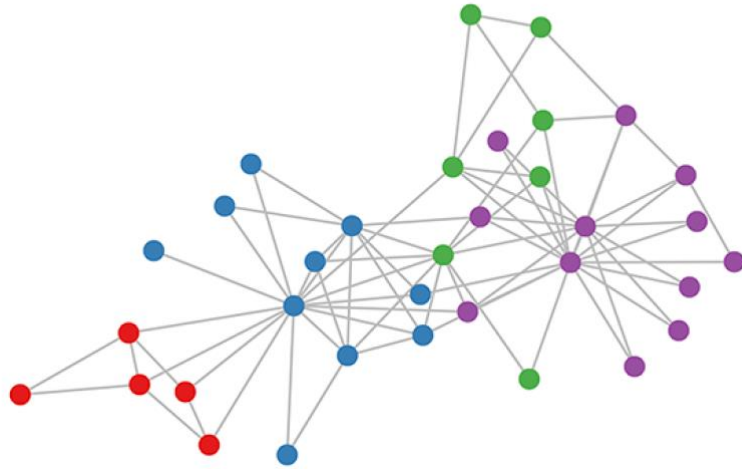


- Goal: Efficient **task-independent feature learning** for machine learning in networks!

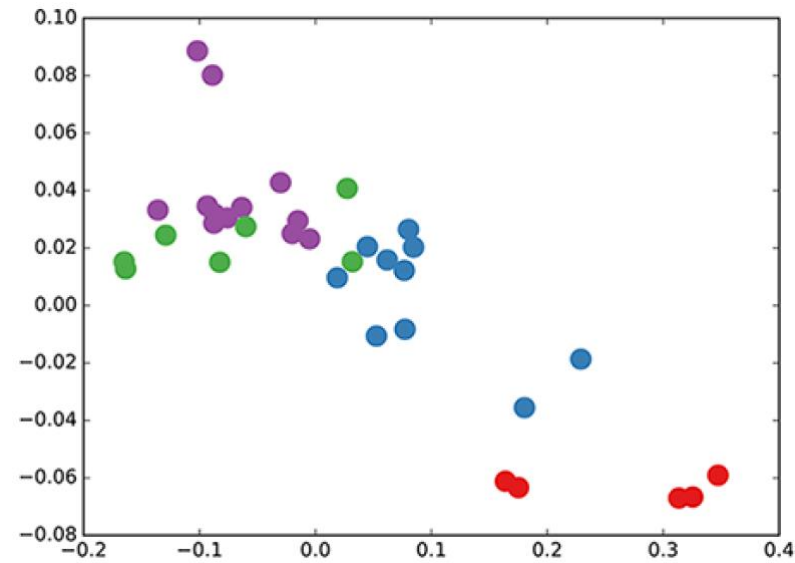




- A good embedding should capture the **graph topology**, **vertex-to-vertex relationship**, and other relevant information about the graph, its subgraphs, and vertices.

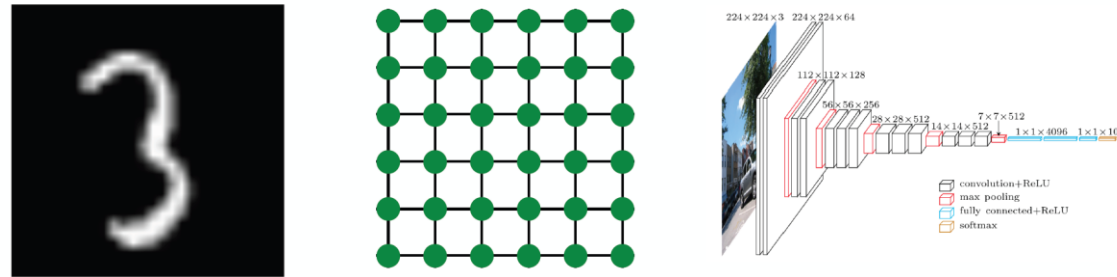


Input

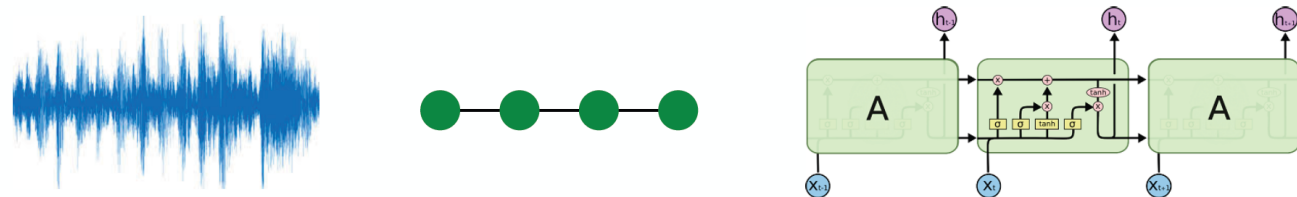


Output

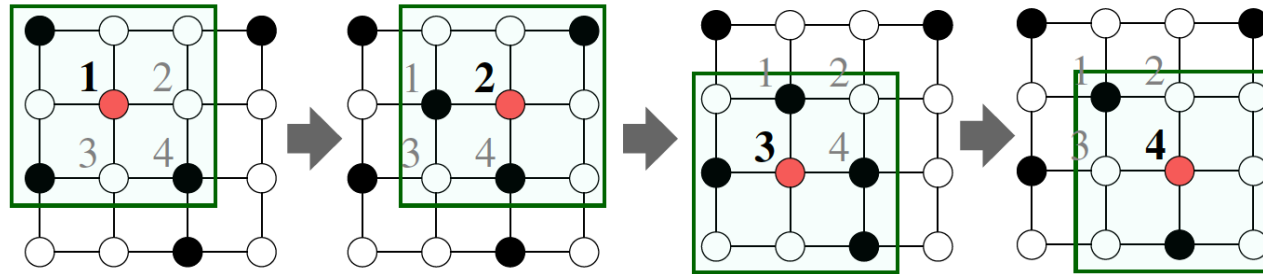
- Modern deep learning toolbox is designed for simple **sequences or grids**.
- CNNs for **fixed-size images/grids**....



- RNNs or word2vec for **text/sequences**...

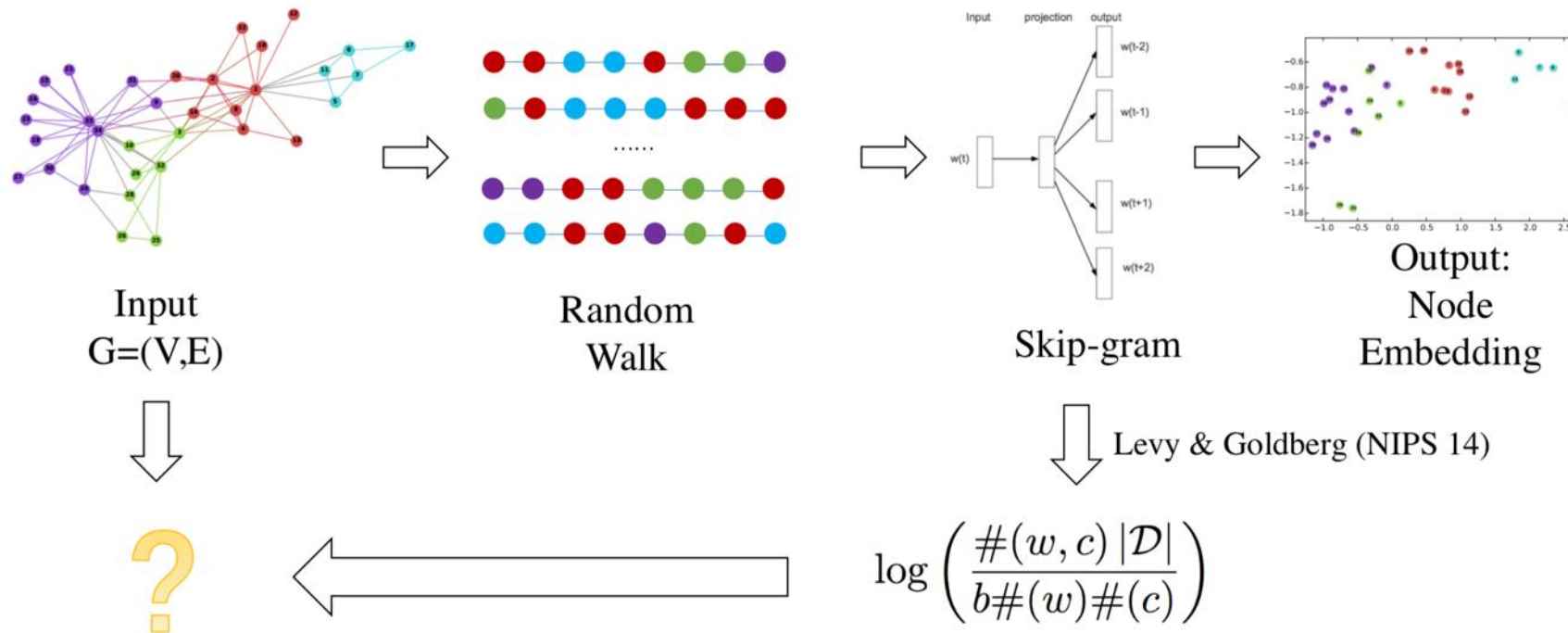


- But networks are far more complex
  - Complex topographical structure (i.e., **no spatial locality** like grids)



- **No fixed node ordering** or reference point (i.e., the isomorphism problem)
- Often dynamic and have multimodal features.

- Random walk approaches
  - Node2Vec
- Multi-hop Similarity
  - LINE
- Deep model (Autoencoder)
  - SDNE
- Subgraph learning
  - Subgraph2Vec



$b$  Number of negative samples

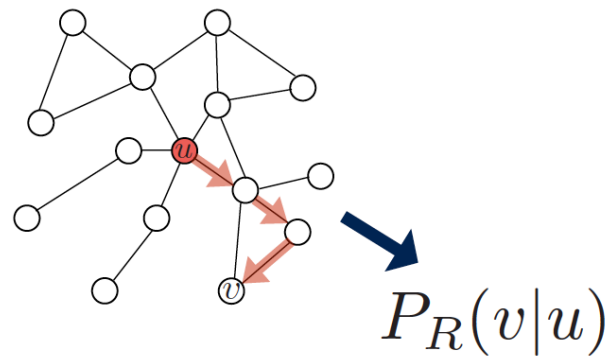
$\#(w, c)$  Co-occurrence of  $w$  and  $c$

$|\mathcal{D}|$  Total number of word-context pairs

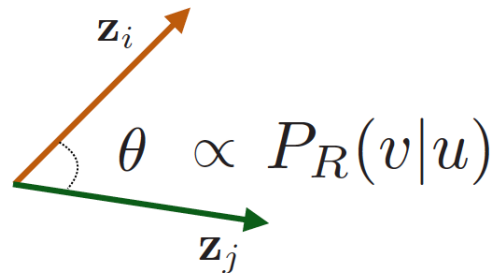
$\#(w)$  Occurrence of word  $w$

$\#(c)$  Occurrence of context  $c$

- Estimate **probability of visiting node  $v$  on a random walk starting from node  $u$**  using some random walk strategy  $R$ .



- Optimize embeddings to encode these random walk statistics.





- Expressivity:
  - Flexible stochastic definition of node similarity that incorporates both local and higher-order neighbourhood information.
- Efficiency:
  - Do not need to consider all node pairs when training; only need to consider pairs that co-occur on random walks.

- Run short **random walks starting from each node** on the graph using some strategy R.
- For each node  $u$  collect  **$N(u)$ , the multiset\* of nodes visited on random walks starting from  $u$ .**
- Optimize embeddings to according to:

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

- \*  $N(u)$  can have repeat elements since nodes can be visited multiple times on random walks.

- Intuition: Optimize embeddings to **maximize likelihood of random walk co-occurrences**.

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

- Parameterize  $P(v | \mathbf{z}_u)$  using **softmax**:

$$P(v|\mathbf{z}_u) = \frac{\exp(\mathbf{z}_u^\top \mathbf{z}_v)}{\sum_{n \in V} \exp(\mathbf{z}_u^\top \mathbf{z}_n)}$$

- Putting things together:
- Optimizing random walk embeddings =  
Finding embeddings  $z_u$  that minimize  $L$

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

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

- But doing this naively is **too expensive**.

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

Nested sum over nodes  
gives  $O(|V|^2)$  complexity!!

The normalization term from  
the softmax is the problem

Can we approximate it?

- Solution: Negative sampling
- i.e., instead of normalizing w.r.t. all nodes, just **normalize against k random “negative samples”**

$$\log \left( \frac{\exp(\mathbf{z}_u^\top \mathbf{z}_v)}{\sum_{n \in V} \exp(\mathbf{z}_u^\top \mathbf{z}_n)} \right)$$
$$\approx \log(\sigma(\mathbf{z}_u^\top \mathbf{z}_v)) - \sum_{i=1}^k \log(\sigma(\mathbf{z}_u^\top \mathbf{z}_{n_i})), n_i \sim P_V$$

sigmoid function

random distribution over all nodes

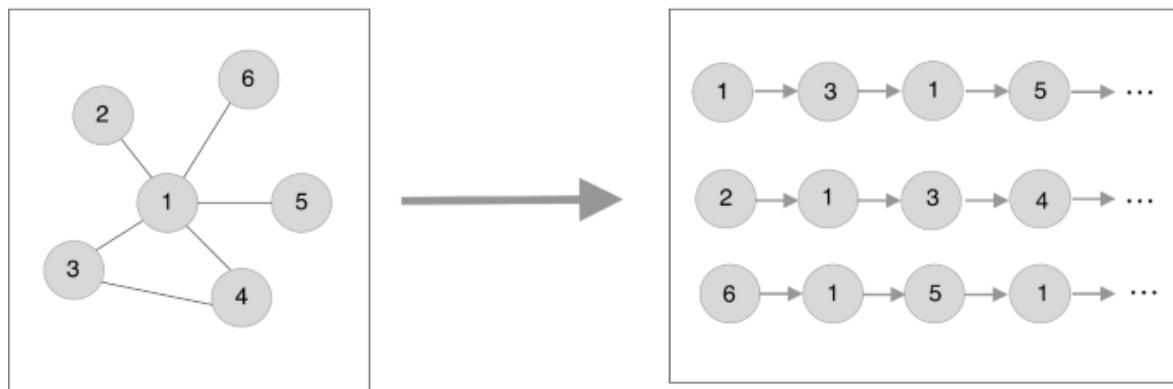


- Run short random walks starting from each node on the graph using some strategy  $R$ .
- For each node  $u$  collect  $NR(u)$ , the multiset of nodes visited on random walks starting from  $u$ .
- Optimize embeddings to according to:

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

We can efficiently approximate this using  
negative sampling!

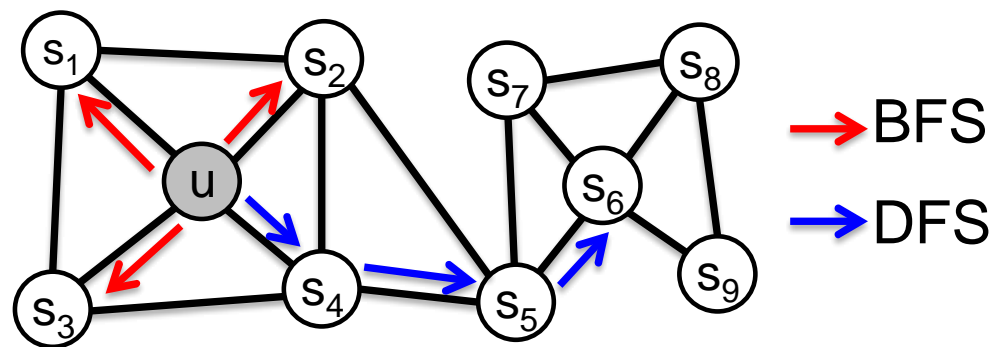
- What strategies should we use to run these random walks?
- Simplest idea: Just **run fixed-length, unbiased random walks** starting from each node (i.e., DeepWalk from Perozzi et al., 2013).



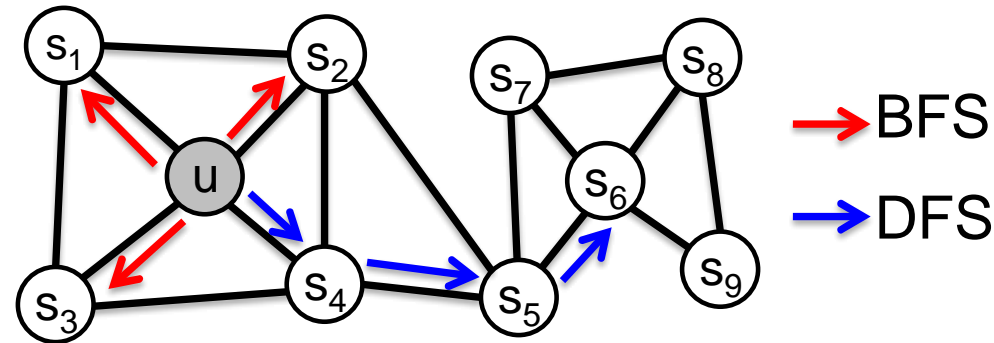
- But can we do better?

➤ Idea:

- use **flexible, biased random walks** that can **trade-off between local and global views** of the network (Grover and Leskovec, 2016).



- Two classic strategies to define a neighborhood  $N(u)$  of a given node  $u$ :

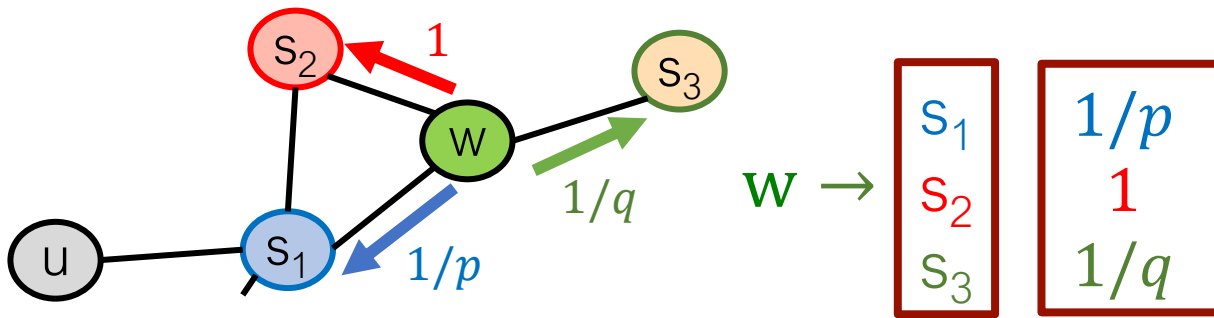


$$N_{BFS}(u) = \{s_1, s_2, s_3\} \quad \text{Local microscopic view}$$

$$N_{DFS}(u) = \{s_4, s_5, s_6\} \quad \text{Global macroscopic view}$$

- Biased random walk  $R$  that given a node  $u$  generates neighborhood  $N(u)$
- Two parameters:
  - Return parameter  $p$ :
    - Return back to the previous node
  - In-out parameter  $q$ :
    - Moving outwards (DFS) vs. inwards (BFS)

- Walker is at  $w$ . Where to go next?



$1/p, 1/q, 1$  are unnormalized probabilities

where:

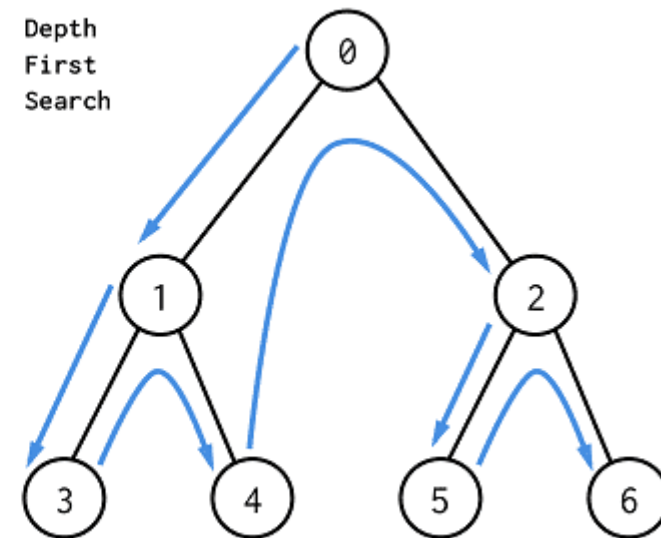
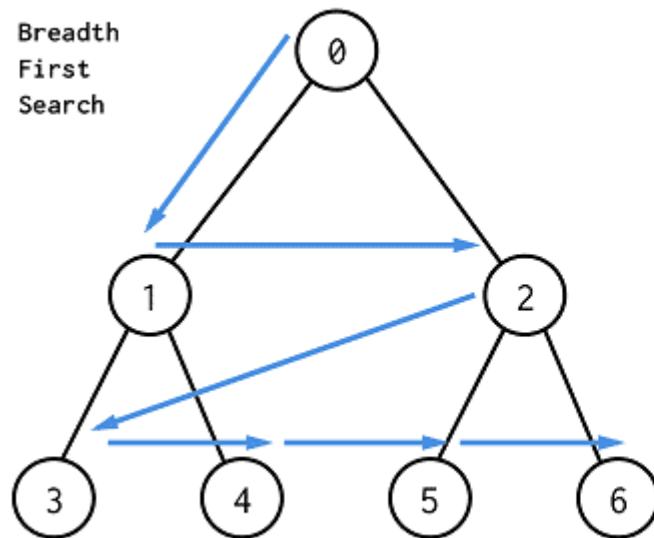
$p, q$ : model transition probabilities

$p$ : return parameter

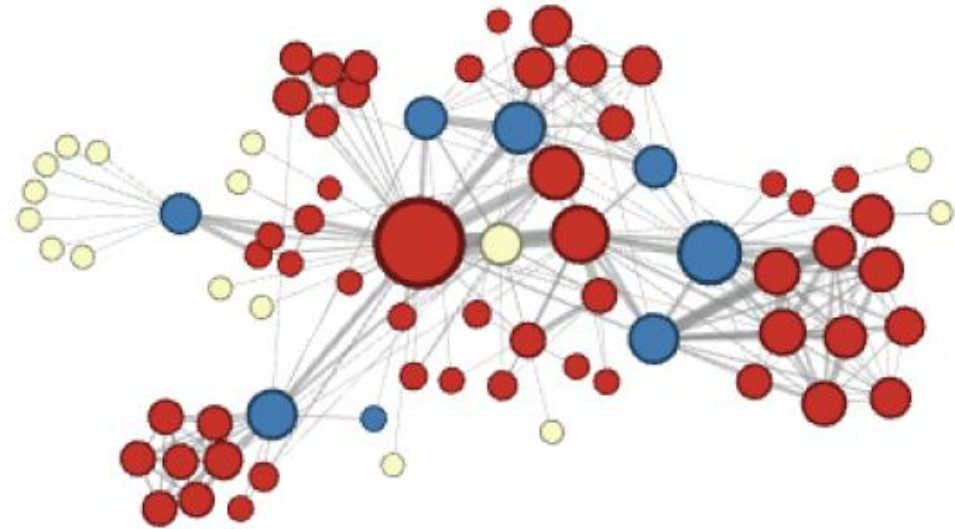
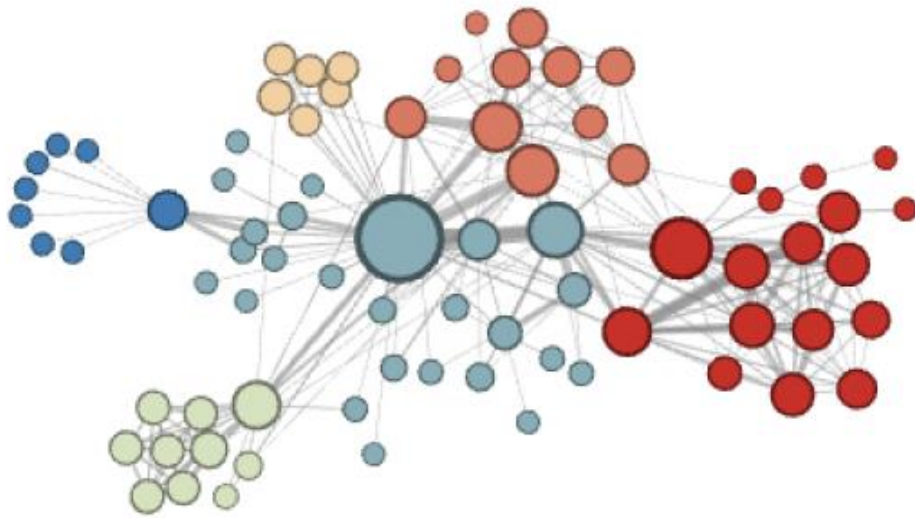
$q$ : "walk away" parameter



- BFS: **Micro-view** of neighbourhood
- DFS: **Macro-view** of neighbourhood

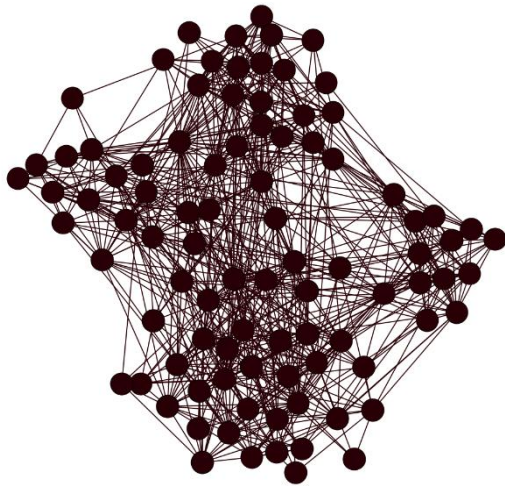


- BFS: **Micro-view** of neighbourhood
- DFS: **Macro-view** of neighbourhood

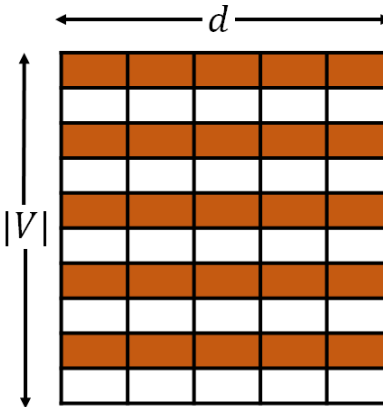


➤ Graph embedding:

**Input Network**



**Node Embeddings**



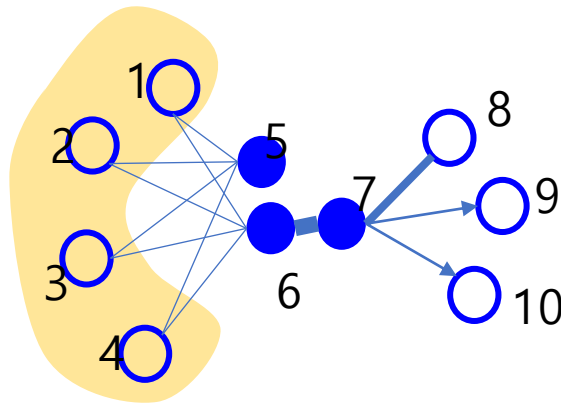
**Data Mining Tasks**

- Classification
- Community Detection
- Link Prediction
- Anomaly Detection
- Sense Making
- ...

Many Possible Applications!

- Has a clear objective function
  - Preserve the **first-order** and **second-order proximity** between the vertices
- Very scalable
  - Effective and efficient optimization algorithm through **asynchronous stochastic gradient descent**
  - Only take a couple of hours to embed network with millions of nodes, billions of edges on a single machine

- The **local pairwise proximity** between the vertices
  - Determined by the **observed links**
- However, many links between the vertices are **missing**
  - **Not sufficient** for preserving the entire network structure

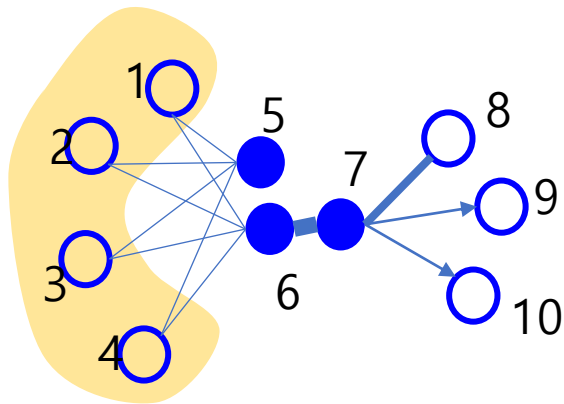


Vertex **6** and **7** have a large first-order proximity

- The **proximity between the neighbourhood structures** of the vertices
- Mathematically, the second-order proximity between each pair of vertices (u,v) is determined by:

$$\hat{p}_u = (w_{u1}, w_{u2}, \dots, w_{u|V|})$$

$$\hat{p}_v = (w_{v1}, w_{v2}, \dots, w_{v|V|})$$



**Vertex 5 and 6 have a large second-order proximity**

$$\hat{p}_5 = (1, 1, 1, 1, 0, 0, 0, 0, 0, 0)$$

$$\hat{p}_6 = (1, 1, 1, 1, 0, 0, 1, 0, 0, 0)$$



- Given an **undirected** edge  $(v_i, v_j)$ , the joint probability of  $v_i, v_j$

$$p_1(v_i, v_j) = \frac{1}{1 + \exp(-\vec{u}_i^T \cdot \vec{u}_j)}$$

$\vec{u}_i$ : Embedding of vertex  $v_i$

$$\hat{p}_1(v_i, v_j) = \frac{w_{ij}}{\sum_{(i', j')} w_{i' j'}}$$

- Objective:

$$O_1 = d(\hat{p}_1(\cdot, \cdot), p_1(\cdot, \cdot))$$

KL-divergence

$$\propto - \sum_{(i, j) \in E} w_{ij} \log p_1(v_i, v_j)$$

- Given a directed edge  $(v_i, v_j)$ , the conditional probability of  $v_j$  given  $v_i$  is:

$$p_2(v_j|v_i) = \frac{\exp(\vec{u}_j'^T \cdot \vec{u}_i)}{\sum_{k=1}^{|V|} \exp(\vec{u}_k'^T \cdot \vec{u}_i)}$$

$\vec{u}_i$ : Embedding of vertex  $i$  when  $i$  is a source node;  
 $\vec{u}_i'$ : Embedding of vertex  $i$  when  $i$  is a target node.

$$\hat{p}_2(v_j|v_i) = \frac{w_{ij}}{\sum_{k \in V} w_{ik}}$$

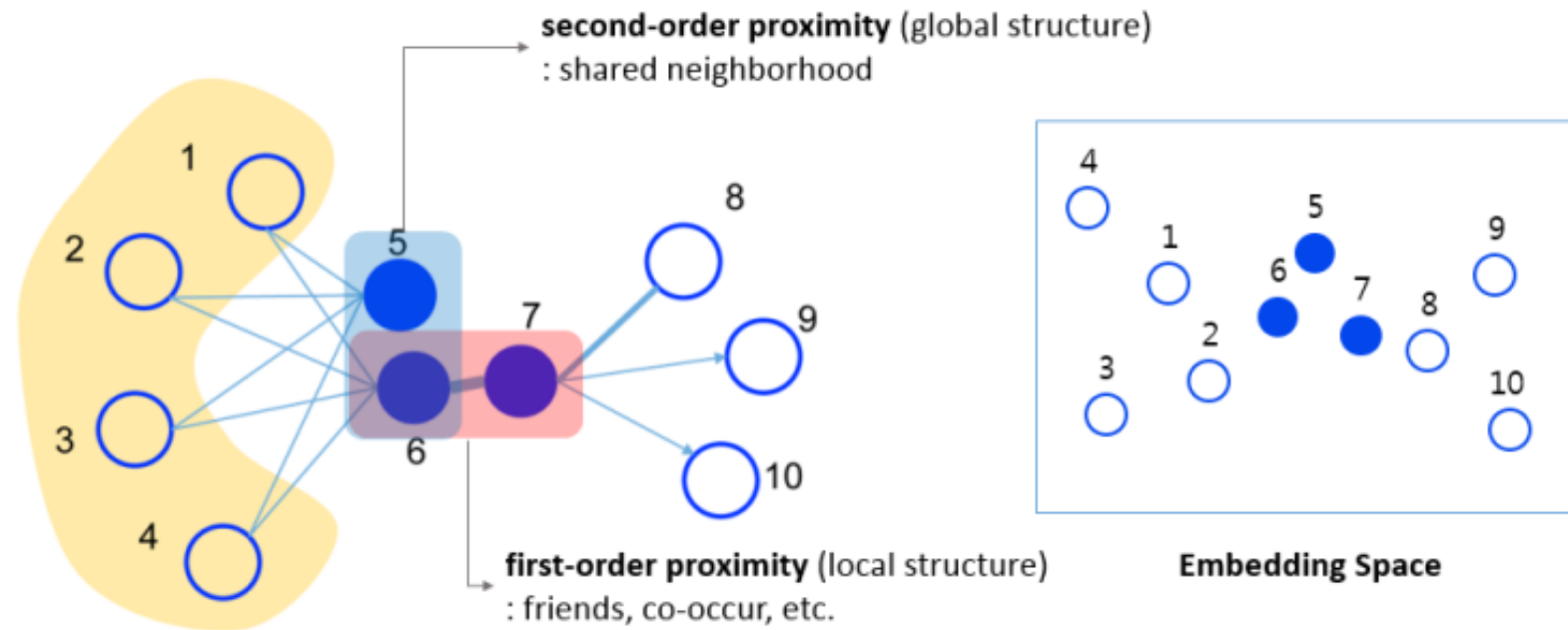
- Objective:

$$O_2 = \sum_{i \in V} \lambda_i d(\hat{p}_2(\cdot | v_i), p_2(\cdot | v_i))$$

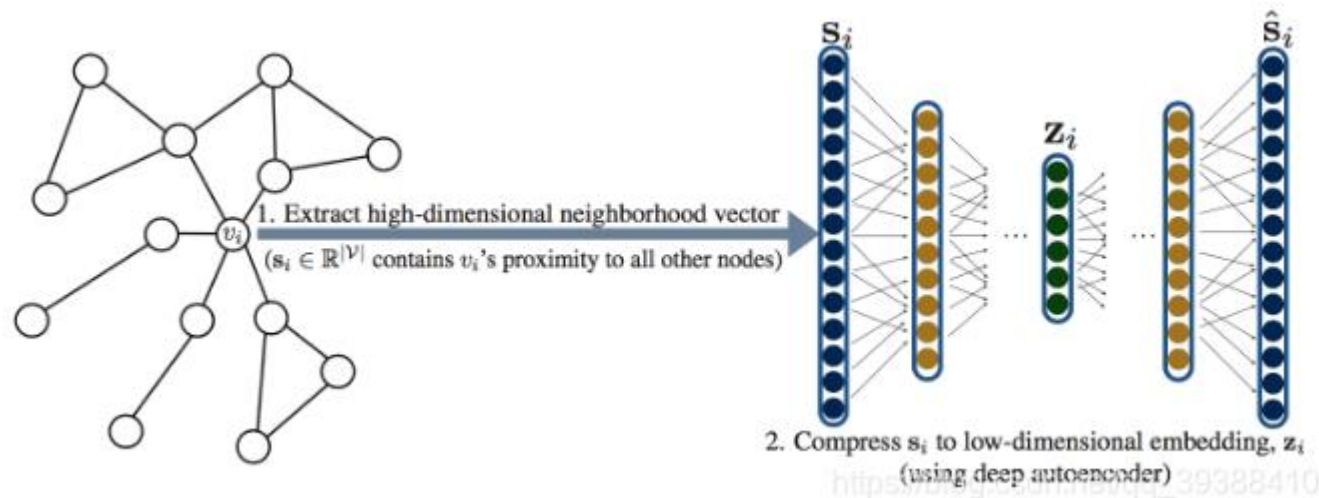
$\lambda_i$ : Prestige of vertex in the network  $\lambda_i = \sum_j w_{ij}$

$$\propto - \sum_{(i,j) \in E} w_{ij} \log p_2(v_j|v_i)$$

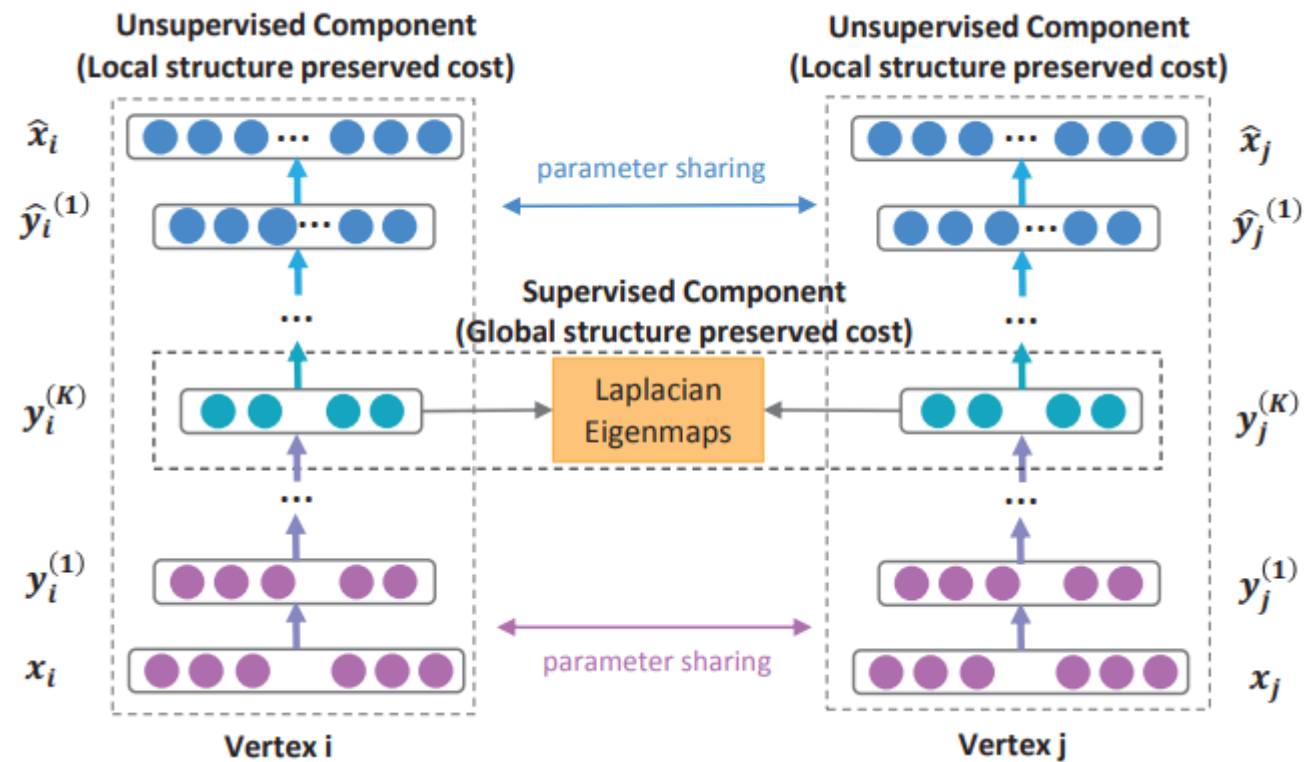
## Preserve both local & global network structure



- Structural Deep Network Embedding based on **Autoencoder**



- The framework of the semi-supervised deep model of SDNE



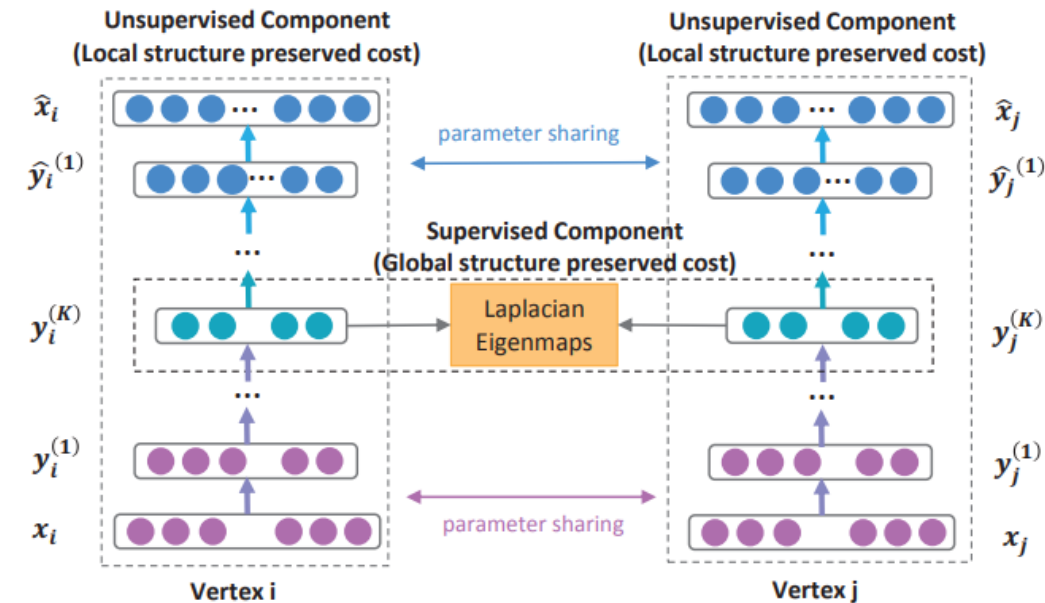
- Then given the input  $x_i$ , the hidden representations for each layer are:

$$y_i^{(1)} = \sigma(W^{(1)}x_i + b^{(1)})$$

$$y_i^{(k)} = \sigma(W^{(k)}y_i^{(k-1)} + b^{(k)}), k = 2, \dots, K$$

- The goal of the autoencoder is to **minimize the reconstruction error of the output and the input.**
- The loss function:

$$\mathcal{L} = \sum_{i=1}^n \|\hat{x}_i - x_i\|_2^2$$

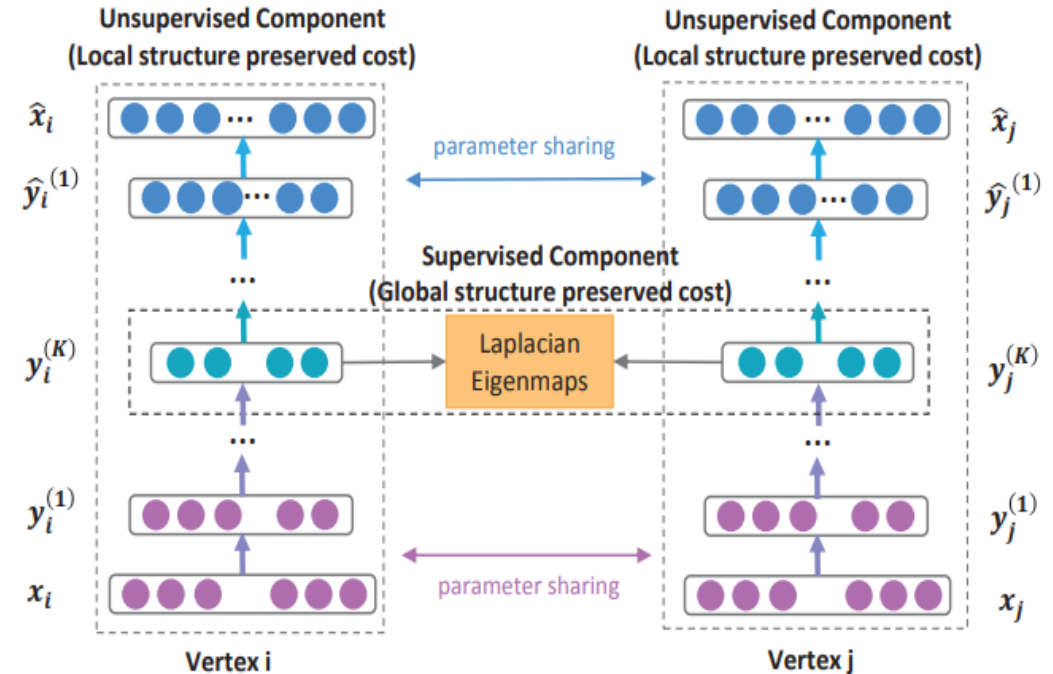


- Loss function for **first-order proximity**:

$$\begin{aligned}\mathcal{L}_{1st} &= \sum_{i,j=1}^n s_{i,j} \|\mathbf{y}_i^{(K)} - \mathbf{y}_j^{(K)}\|_2^2 \\ &= \sum_{i,j=1}^n s_{i,j} \|\mathbf{y}_i - \mathbf{y}_j\|_2^2\end{aligned}$$

- Impose more penalty to the reconstruction **error of the non-zero elements** than that of zero elements:

$$\begin{aligned}\mathcal{L}_{2nd} &= \sum_{i=1}^n \|(\hat{\mathbf{x}}_i - \mathbf{x}_i) \odot \mathbf{b}_i\|_2^2 \\ &= \|(\hat{X} - X) \odot B\|_F^2\end{aligned}$$

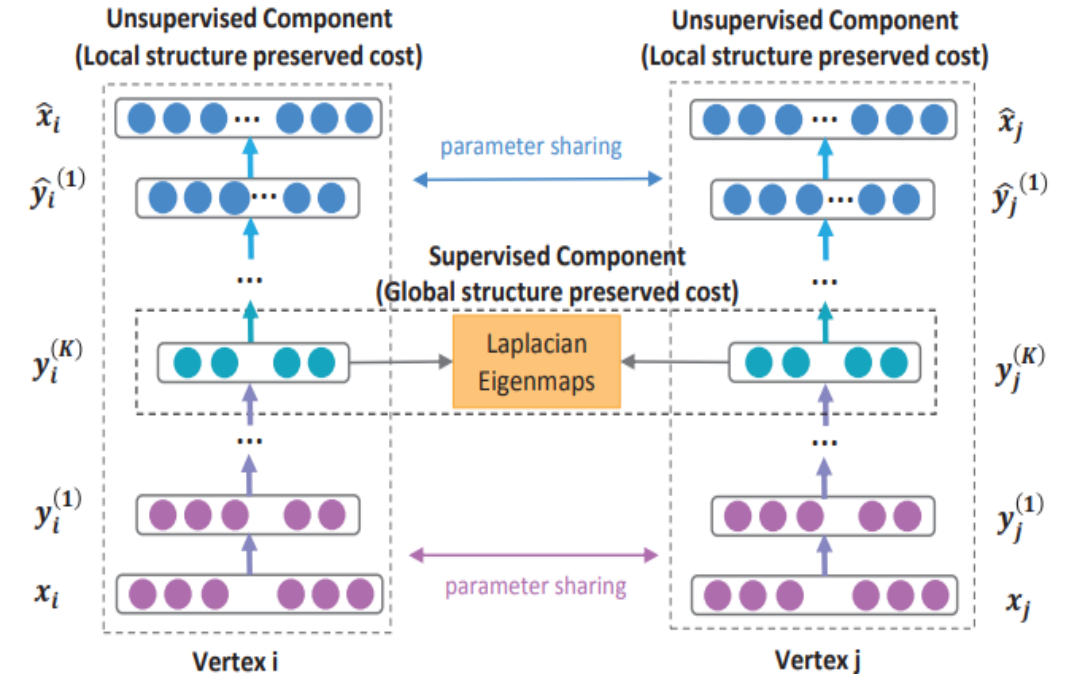


- To preserve the **first-order and second-order proximity** simultaneously, we need to minimize the **joint loss**:

$$\begin{aligned}\mathcal{L}_{mix} &= \mathcal{L}_{2nd} + \alpha \mathcal{L}_{1st} + \nu \mathcal{L}_{reg} \\ &= \|(\hat{X} - X) \odot B\|_F^2 + \alpha \sum_{i,j=1}^n s_{i,j} \|\mathbf{y}_i - \mathbf{y}_j\|_2^2 + \nu \mathcal{L}_{reg}\end{aligned}$$

- where  $\mathcal{L}_{reg}$  is an L2-norm **regularizer** term to prevent overfitting, which is defined as follows:

$$\mathcal{L}_{reg} = \frac{1}{2} \sum_{k=1}^K (\|W^{(k)}\|_F^2 + \|\hat{W}^{(k)}\|_F^2)$$





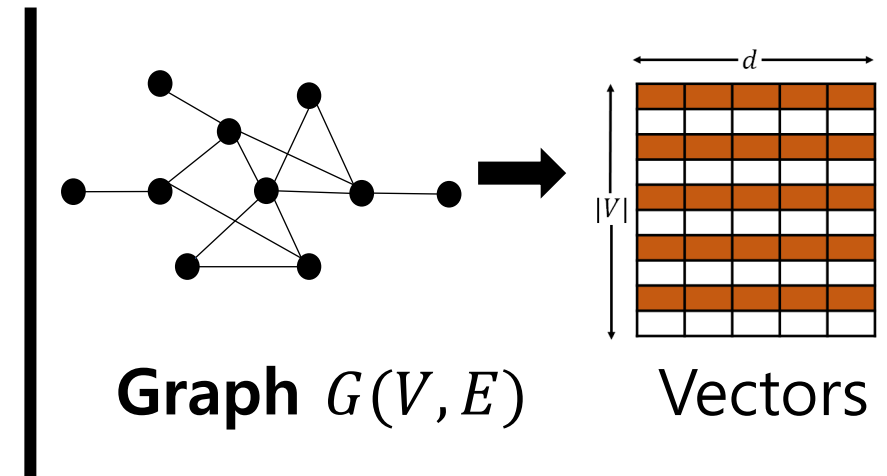
➤ Popular previous works include

DeepWalk[Perozzi+, KDD2014]

Node2vec[Grover+, KDD 2016]

SDNE[Wang+, KDD 2016]

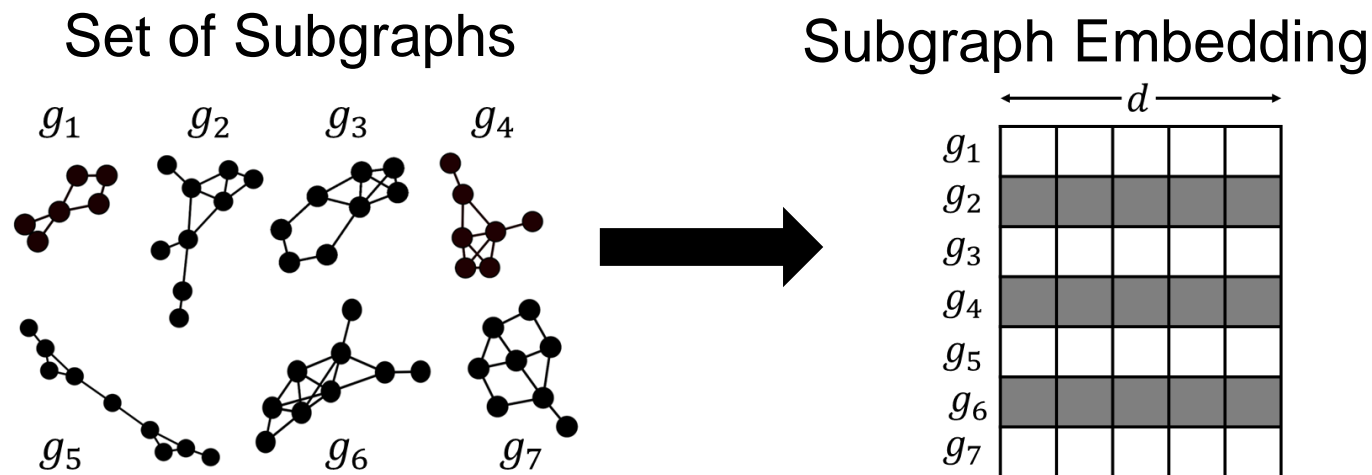
LINE[Tang+, WWW 2015]



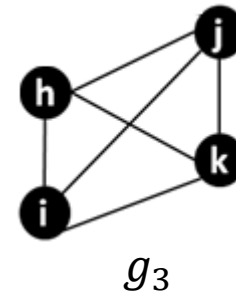
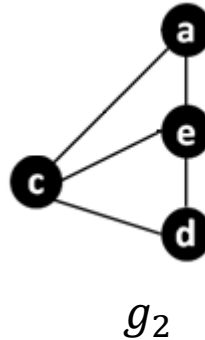
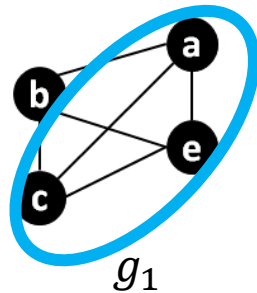
Limited just to the node embeddings

- Learning **representation of substructures**
  - Extend the **WL relabeling strategy** to define a **proper context** for a given subgraph.
  - A modification to the skipgram model enabling it to **capture varying length radial contexts**

- Given
  - A set  $S = \{g_1, g_2, \dots, g_n\}$  of subgraphs
  - Typically for the same graph
  - An integer  $d$
- Learn
  - $d$ -dimensional embedding for each subgraph
  - Such that **pre-defined subgraph property is preserved**



- What subgraph property to preserve?
  - Neighbourhood Property:
    - Captures **neighbourhood information within the subgraph**



- Subgraph  $g_1$  and  $g_2$  share neighbourhood
- Subgraph  $g_3$  does not

- Generate rooted subgraphs around every node in a given graph
- Considers all the rooted subgraphs (up to a certain degree) of neighbours of  $r$  as the context of target subgraphs

---

**Algorithm 2:** GETWLSUBGRAPH ( $v, G, d$ )

---

**input** :  $v$ : Node which is the root of the subgraph  
 $G = (V, E, \lambda)$ : Graph from which subgraph has to be extracted  
 $d$ : Degree of neighbours to be considered for extracting subgraph

**output:**  $sg_v^{(d)}$ : rooted subgraph of degree  $d$  around node  $v$

```
1 begin
2    $sg_v^{(d)} = \{\}$ 
3   if  $d = 0$  then
4      $sg_v^{(d)} := \lambda(v)$ 
5   else
6      $\mathcal{N}_v := \{v' \mid (v, v') \in E\}$ 
7      $M_v^{(d)} := \{\text{GETWLSUBGRAPH}(v', G, d - 1) \mid v' \in \mathcal{N}_v\}$ 
8      $sg_v^{(d)} := sg_v^{(d)} \cup \text{GETWLSUBGRAPH}$ 
        $(v, G, d - 1) \oplus \text{sort}(M_v^{(d)})$ 
9   return  $sg_v^{(d)}$ 
```

---

- The skipgram model maximizes **co-occurrence probability among the subgraphs** that appear within a given context window.

---

## Algorithm 3: RADIALSKIPGRAM ( $\Phi, sg_v^{(d)}, G, D$ )

---

```

1 begin
2    $context_v^{(d)} = \{\}$ 
3   for  $v' \in \text{NEIGHBOURS}(G, v)$  do
4     for  $\partial \in \{d-1, d, d+1\}$  do
5       if  $(\partial \geq 0 \text{ and } \partial \leq D)$  then
6          $context_v^{(d)} = context_v^{(d)} \cup$ 
           GETWLSUBGRAPH( $v', G, \partial$ )
7   for each  $sg_{cont} \in context_v^{(d)}$  do
8      $J(\Phi) = -\log \Pr (sg_{cont} | \Phi(sg_v^{(d)}))$ 
9      $\Phi = \Phi - \alpha \frac{\partial J}{\partial \Phi}$ 

```

---





네트워크 과학연구실  
NETWORK SCIENCE LAB



가톨릭대학교  
THE CATHOLIC UNIVERSITY OF KOREA

