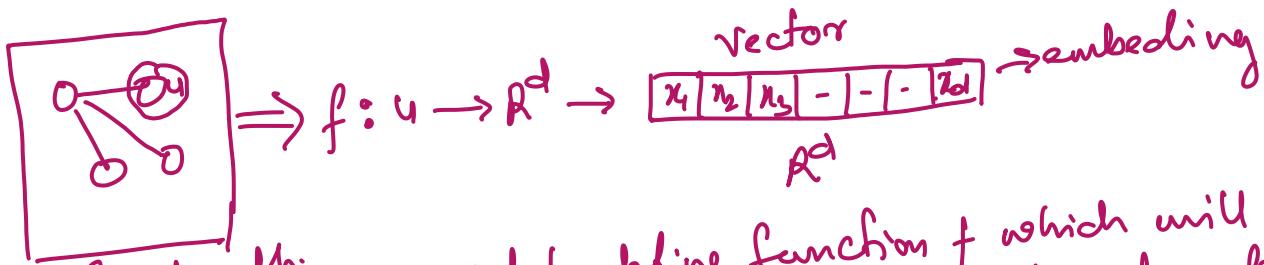


## Node Embedding

Till now we have seen hand-crafted method to generate feature set for nodes, links, Graph where main idea was that based on structural property two features should be closer if structural property is same. But we have seen constructing hand-crafted feature set is computationally expensive and time consuming. To overcome this we want mechanism which will automatically extract the features which is relevant for the graph. Hand Crafted feature extraction is called Feature engineering. When we extract features automatically then that is called Graph representation learning.



Input-graph: we need to define function  $f$  which will read input-node  $u$  and transform it in  $d$ -dimensional real number vector space. This vector can be used for various downstream tasks like

- ① Node classification
- ② Link prediction
- ③ Graph classification
- ④ Anomalous node detection
- ⑤ Clustering etc.

similarity in embedding also indicate similarity of node in network too. Process of embedding is also called encoding network information.

Node Embedding framework:

3 component:

- ① Encoder
- ② Node similarity function
- ③ Decoder

① Encoder: Process of generating  $R^d$  vector space embedding.

② In original network we need to define node similarity function which measures the similarity in original network.

(3) Decoder: Maps the embedding to similarity score.

(4) Optimize the parameter of encoder so that-

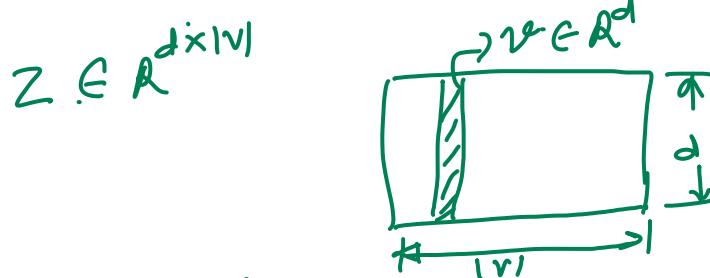
similarity( $u, v$ )  $\approx z_v^T z_u \rightarrow$  similarity in embedding space.

$ENC(u) \rightarrow z_u$   $\downarrow$   
similarity in original  
networks

We will study two kind of node embedding:

- (1) Random walk
- (2) Node2Vec

when size of embedding vector is  $d$  and suppose we have  $|V|$  number of nodes then all will be stacked in single matrix  $Z$



Hence each column represent one nodes embedding vector.

Since embedding is lookups in Matrix  $Z$  so this is called shallow embedding. When size of graph is large (billions of node) then size of  $Z$  becomes issue. Each element of  $Z$  is parameter that needs to learn.

Here objective function is to maximize  $z_v^T z_u$  as per similarity function for actual network. So decoder is designed based on similarity function.

Similarity function can be defined based on many

criterion e.g

Nodes are linked

share neighbours

have similar structural role etc.

Here node embedding is task independent function and we do not use node label, node feature etc. Our aim is to create node embedding which captures some aspect of networks structure by encoder. This embedding can be used for any task.

## (Deepwalk) Random Walk Approaches for Node Embeddings

- ① Random walk strategy: We need to fix how randomly we are going to walk over graph. One simple way can be pick randomly some node and choose randomly some next node and keep on walking till predefined fixed length not finished.



- ② We need to randomly initialize  $z_u$  vector of length  $d$  with real value which will be optimized based on the idea that similar node embedding will have closer value.

- ③ Probability of reaching node  $v$  from  $u$  using random walk strategy & will be multiplication of each step as all hop is independent - because its random walk.

$$p(v|z_u) = p(v_j|z_u) \times p(v_k|z_{v_j}) \times p(v|z_{v_k}) \rightarrow \text{similarity func.}$$

if we take log likelihood then we can generalize as

$$= \sum \log p(N_{\Delta}(u)|z_u)$$

- ④ In random walk if probability of some node occurrence is high means they are closer and consider as similar.

- ⑤ Using dot product - we can find similarity score
- $$z_u^T z_v \xrightarrow{\text{to}} z_v^T P(v|z_u)$$
- We can use cosθ as encoding function.

- ⑥ Our objective is to learn  $d$ -dimensional embedding vector that preserves the similarity.  
Here if nodes are nearby then their embedding should be similar.

- ⑦ Let  $G(V, E)$  is graph where we need to define
- $$f: u \rightarrow \mathbb{R}^d, f(u) = z_u$$

- ⑧ Since nearby nodes should have similar embedding then we need to maximize log-likelihood function
- $$\max \sum_{u \in V} \log p(N_{\Delta}(u)|z_u)$$
- where  $N_{\Delta}(u)$  is neighbour nodes as per strategy  $\Delta$ .

Here  $\text{NA}(u)$  will be multiset - because of strategy A.

- ⑨ Now our problem converted into optimization problem which can be written as minimization problem.

$$f = \sum_{u \in V} \sum_{v \in \text{NA}(u)} -\log(P(v|z_u)) \rightarrow \text{for given fixed length.}$$

For all node  $v \in G(V, E) \Rightarrow u \in V$

- ⑩ Similarity probability can be defined using softmax function.

$z_u^T z_v \Rightarrow$  Gives similarity score.

$$\text{Softmax function} \Rightarrow \frac{\exp(z_u^T z_v)}{\sum_{n \in V} \exp(z_u^T z_n)} \Rightarrow \text{This must be same as } P_a(v|z_u)$$

$$\Rightarrow P_a(v|z_u) = \frac{\exp(z_u^T z_v)}{\sum_{n \in V} \exp(z_u^T z_n)} \rightarrow \text{node embedding similarity function}$$

similarity function  
of actual network

- ⑪ Now optimization function be written as

$$f = \sum_{u \in V} \sum_{v \in \text{NA}(u)} -\log \frac{\exp(z_u^T z_v)}{\sum_{n \in V} \exp(z_u^T z_n)} \Rightarrow \text{Minimization problem.}$$

- ⑫ Solving this naively will be very expensive because nesting of summation over all nodes. Complexity  $O(|V|^2)$

- ⑬ Normalization in softmax  $\sum_{n \in V} \exp(z_u^T z_n)$  is culprit. we need to approximate this. This can be done using negative sampling. We do not take all nodes but  $k$ -nodes for approximation and instead of exponential function we use sigmoid function ( $\sigma$ ) for probability.

$$\log \frac{\exp(z_u^T z_v)}{\sum_{n \in V} \exp(z_u^T z_n)} \approx \log \left[ \frac{\sigma(z_u^T z_v)}{\sum_{n \in K} \sigma(z_u^T z_n)} \right]$$

$$= \log \sigma(z_u^T z_v) - \sum_{i=1}^k \log (\sigma(z_u^T z_{n_i}))$$

Here  $\eta_i$  is random distribution over nodes.

By using  $k$ -sample improves the efficiency of calculation.

④ Higher  $k$  gives more robust estimate

⑤ Higher  $k$  corresponds to higher bias on negative ~~est~~ events.

⑥ Sample  $K$  negative nodes each with probability proportional to its degree.

⑦  $K=5-20$  suggested

⑧ Usually negative sampling for any node alone but most efficient way is to take nodes which are not on the walk.

14 Once optimization function defined then using gradient descent- we can optimize  $L$  w.r.t  $z_u$

$$\frac{\partial L}{\partial z_u}$$

$$z_{ut} \leftarrow z_{ut} - \eta \frac{\partial L}{\partial z_u} \quad \text{when } \eta \text{ is learning rate.}$$

15 By stochastic gradient-descent- method we can train  $z_u$  by optimizing  $L$  for given training examples.

We choose random walk because this is expressive due to giving higher probability to closer nodes. Other reason is we do not need to consider all pair node but only those node which is participating in random walk.

## NODE2VEC

Now question came, can we tune the random walk more so that embedding is more expressive? How to make random walk (DeepWalk) more richer? currently using uniformly distributed random walk which may be too constraint- for downstream task.



























