



中國人民大學  
RENMIN UNIVERSITY OF CHINA



# Network Embedding: Recent Progress and Applications

Xin Zhao

[batmanfly@qq.com](mailto:batmanfly@qq.com)

Renmin University of China

# Outline

- Preliminaries
  - word2vec
- Network Embedding Models
  - DeepWalk
  - Node2vec
  - GENE
  - LINE
  - SDNE
- Applications of Network Embedding
  - Basic applications
  - Visualization
  - Text classification
  - Recommendation
- Conclusion

# Preliminaries

- Softmax functions
- Distributional semantics
- Word2vec
  - CBOW
  - Skip-gram

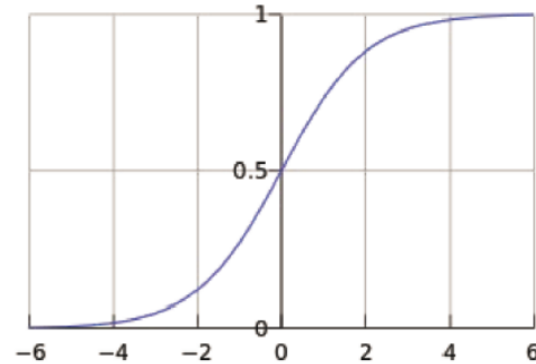
# Preliminaries

- Representation learning
  - Using machine learning techniques to derive data representation
- Distributed representation
  - Different from one-hot representation, it uses dense vectors to represent data points
- Embedding
  - Mapping information entities into a low-dimensional space

# Sigmoid Function

- A sigmoid function can map a real value to the interval (0, 1)

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$



# Softmax function

- It transforms a  $K$ -dimensional real vector into a *probability distribution*
  - A common transformation function to derive objective functions for classification or discrete variable modeling

$$\sigma(\mathbf{z})_j = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}} \quad \text{for } j = 1, \dots, K$$

# Distributional semantics

- Target word = “stars”

he curtains open and the stars shining in on the barely  
ars and the cold , close stars " . And neither of the w  
rough the night with the stars shining so brightly , it  
made in the light of the stars . It all boils down , wr  
surely under the bright stars , thrilled by ice-white  
sun , the seasons of the stars ? Home , alone , Jay pla  
m is dazzling snow , the stars have risen full and cold  
un and the temple of the stars , driving out of the hug  
in the dark and now the stars rise , full and amber a  
bird on the shape of the stars over the trees in front  
But I could n't see the stars or the moon , only the  
they love the sun , the stars and the stars . None of  
r the light of the shiny stars . The plash of flowing w  
man 's first look at the stars ; various exhibits , aer  
rief information on both stars and constellations, inc

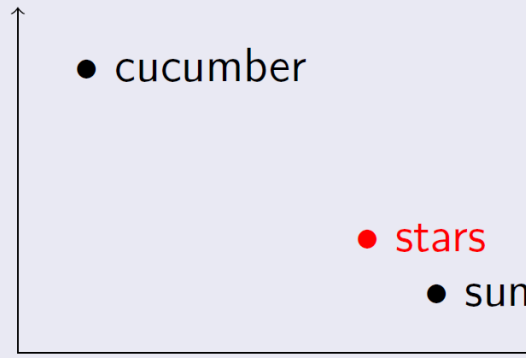
# Distributional semantics

- Collect the contextual words for “stars”

## Construct vector representations

	shining	bright	trees	dark	look
stars	38	45	2	27	12

## Similarity in meaning as vector similarity





# Word2Vec

- Input: a sequence of words from a vocabulary  $V$
- Output: a fixed-length vector for each term in the vocabulary
  - $\mathbf{v}_w$

It implements the idea of distributional semantics using a shallow neural network model.

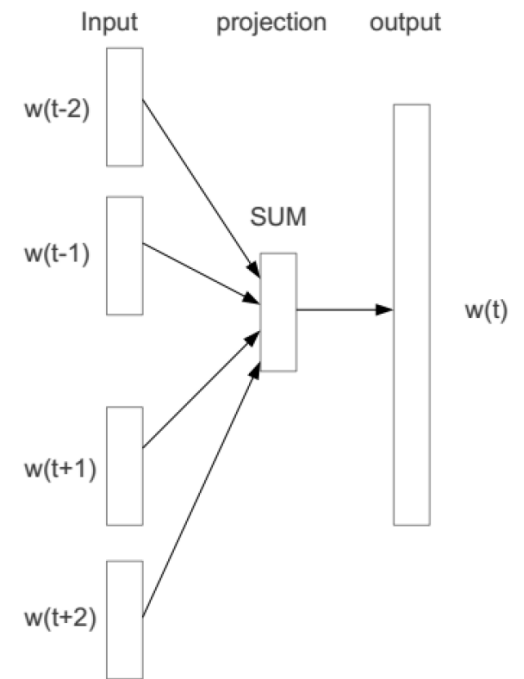
# Architecture 1: CBOW

- **CBOW** predicts the current word using surrounding contexts

–  $Pr(w_t | \text{context}(w_t))$

- Window size  $2c$

- $\text{context}(w_t) = [w_{t-c}, \dots, w_{t+c}]$



# Architecture 1: CBOW

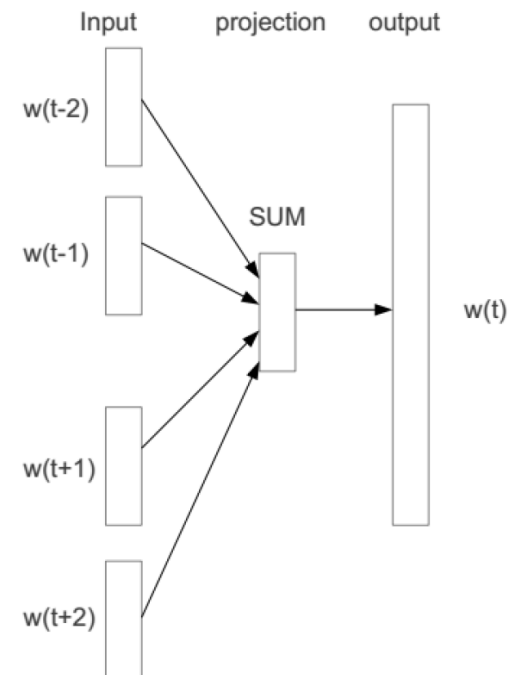
- **CBOW** predicts the current word using surrounding contexts

- $Pr(w_t | \text{context}(w_t))$

- Using a  $K$ -dimensional vector to represent words

- $w_t \rightarrow v_{w_t}$

- $\tilde{v}_{w_t} = \frac{\sum_{i=t-c}^{t+c} v_{w_i}}{2c} \quad (i \neq t)$



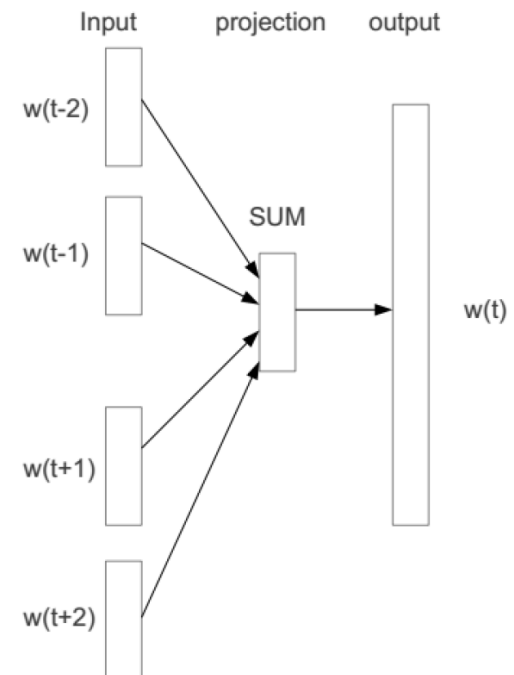
# Architecture 1: CBOW

- **CBOW** predicts the current word using surrounding contexts

- $Pr(w_t | \text{context}(w_t))$

- Basic Idea

- Given the context of the current word  $\tilde{v}_{w_t}$
    - $\text{Sim}(\tilde{v}_{w_t}, v_{w_t}) > \text{Sim}(\tilde{v}_{w_t}, v_{w_j})$



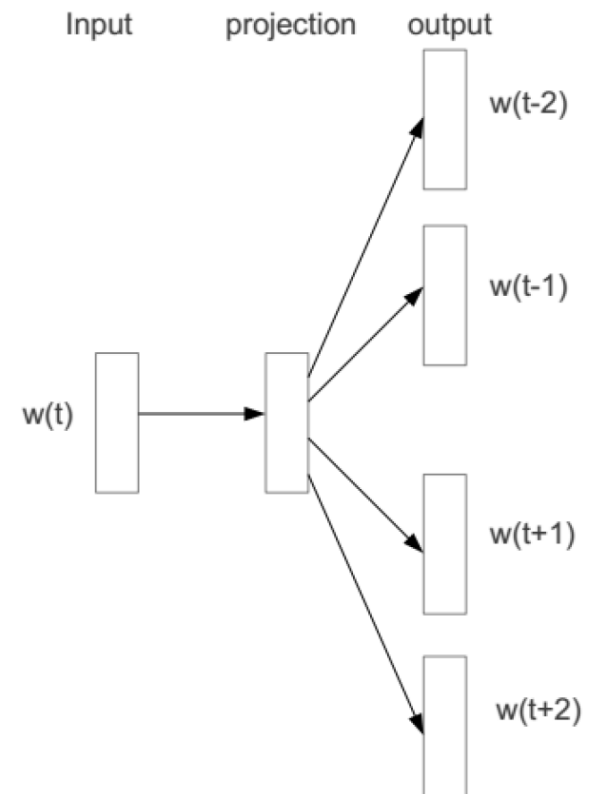
# Architecture 1: CBOW

- How to formulate the idea
  - Using a softmax function
  - Considered as a classification problem
    - Each word is a classification label

$$P(w|w_{\text{context}}) = \frac{\exp(\text{sim}(\tilde{v}_w, v_w))}{\sum_{w'} \exp(\text{sim}(\tilde{v}_w, v_{w'}))}$$

# Architecture 2

- **Skip-gram** predicts surrounding words using the current word
  - $Pr(\text{context}(w_t) \mid w_t)$ 
    - Window size  $2c$
    - $\text{context}(w_t) = [w_{t-c}, \dots, w_{t+c}]$



# Architecture 2

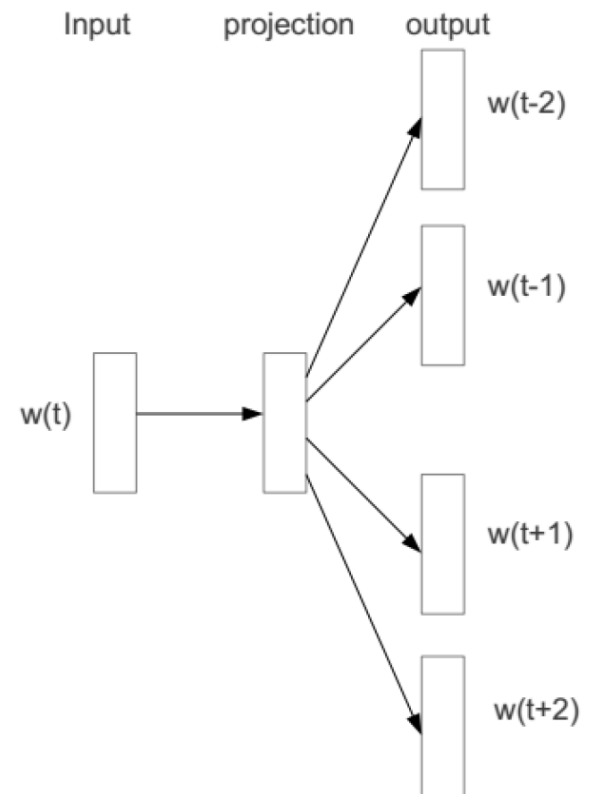
- **Skip-gram** predicts surrounding words using the current word

–  $Pr(\text{context}(w_t) \mid w_t)$

- Window size  $2c$

- $\text{context}(w_t) = [w_{t-c}, \dots, w_{t+c}]$

$$P(w'|w) = \frac{\exp(\text{sim}(\mathbf{v}_w, \mathbf{v}_{w'}))}{\sum_{w'} \exp(\text{sim}(\mathbf{v}_w, \mathbf{v}_{w'}))}$$



# Time Complexity for Direct Optimization

## CBOW

$$P(w|w_{\text{context}}) = \frac{\exp(\text{sim}(\tilde{\mathbf{v}}_w, \mathbf{v}_w))}{\sum_{w'} \exp(\text{sim}(\tilde{\mathbf{v}}_w, \mathbf{v}_{w'}))}$$

## Skip-Gram

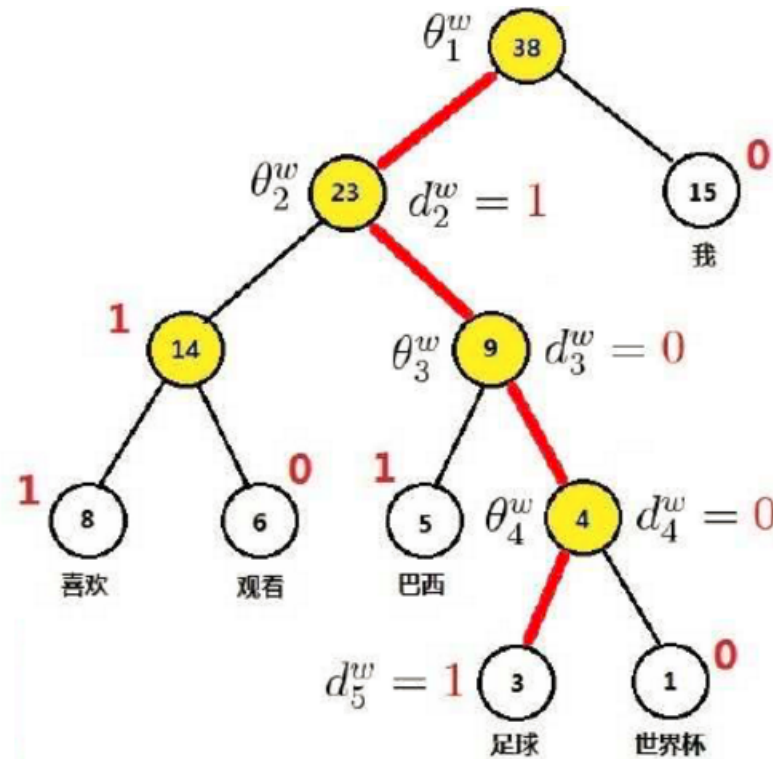
$$P(w'|w) = \frac{\exp(\text{sim}(\mathbf{v}_w, \mathbf{v}_{w'}))}{\sum_{w''} \exp(\text{sim}(\mathbf{v}_w, \mathbf{v}_{w''}))}$$

$$O(|V|)$$



# Optimization I: Hierarchical Softmax

$O(\log_2 |V|)$



It can be verified the sum of the probabilities for all the term nodes is equal to 1

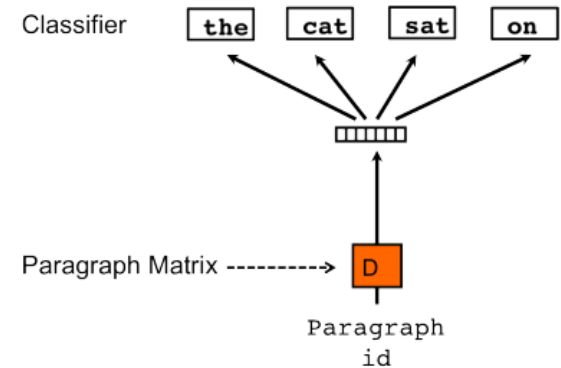
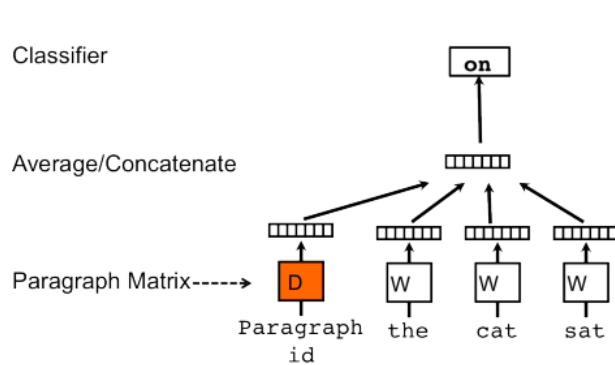
# Optimization II: Negative Sampling

$O(1+k)$

$$\log \sigma(v'_{w_O}{}^\top v_{w_I}) + \sum_{i=1}^k \mathbb{E}_{w_i \sim P_n(w)} \left[ \log \sigma(-v'_{w_i}{}^\top v_{w_I}) \right]$$

# Doc2Vec

- Two architectures



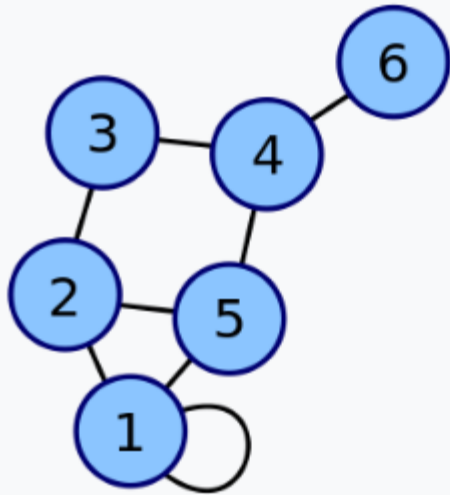
– Both the representations for docs and words can be obtained

# Network Embedding Models

- DeepWalk
- Node2vec
- GENE
- LINE
- SDNE

# Traditional Network Representation

- How to present a graph

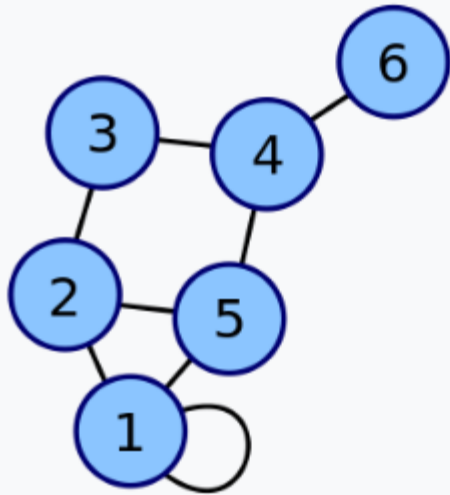


$$\begin{pmatrix} 2 & 1 & 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \end{pmatrix}$$

Coordinates are 1 - 6.

# Traditional Network Representation

- How to present a graph

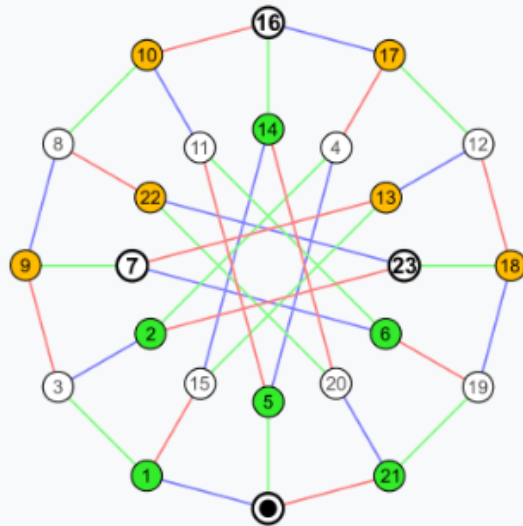


$$\begin{pmatrix} 2 & 1 & 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \end{pmatrix}$$

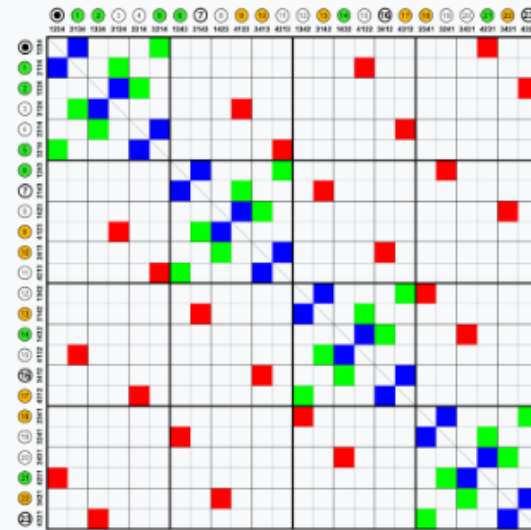
Coordinates are 1 - 6.

# Traditional Network Representation

- How to present a graph



Nauru graph



Coordinates are 0 - 23.

White fields are zeros, colored fields are ones.

# Traditional Network Representation

- Classical graph embedding algorithms
  - MDS, IsoMap, LLE, Laplacian Eigenmap
  - **Hard to scale up**

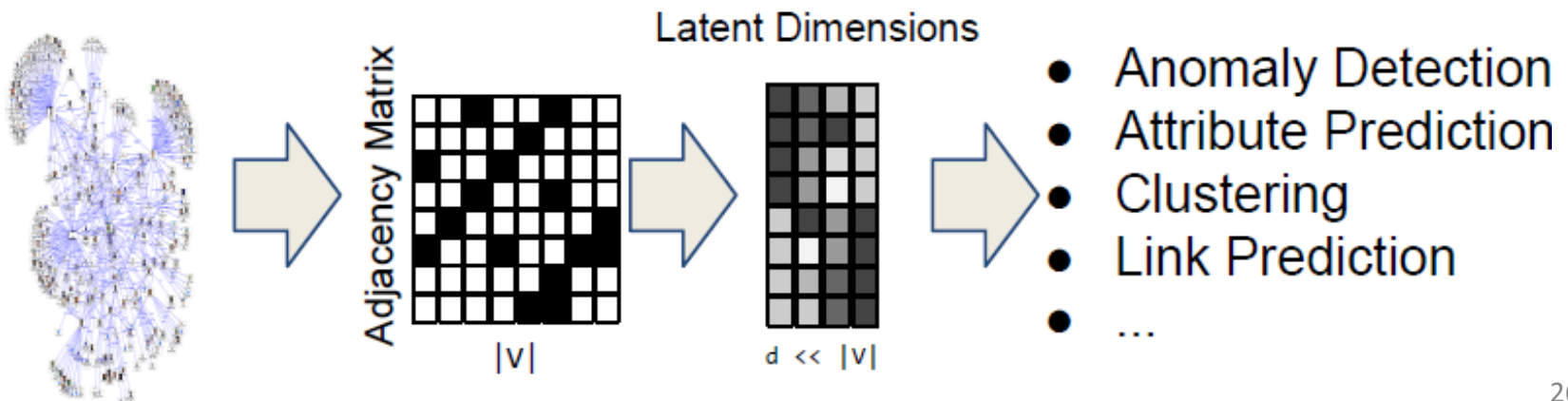


# Network Embedding Models

- **DeepWalk** (Perozzi et al., KDD 2014)
- Node2vec
- GENE
- LINE
- SDNE

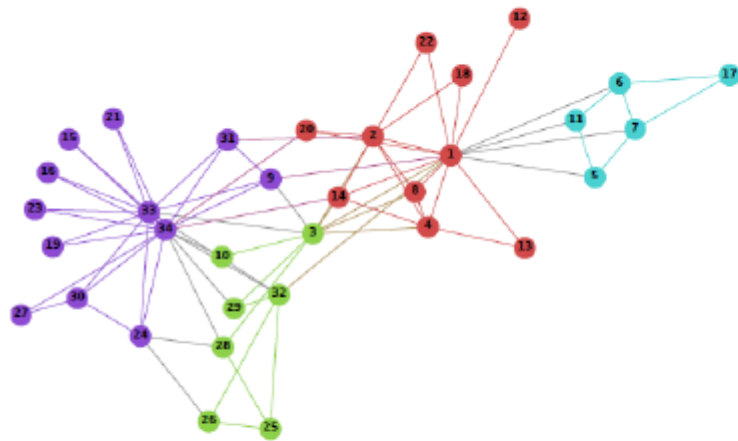
# What is network embedding?

- We map each node in a network into a low-dimensional space
  - Distributed representation for nodes
  - Similarity between nodes indicate the link strength
  - Encode network information and generate node representation

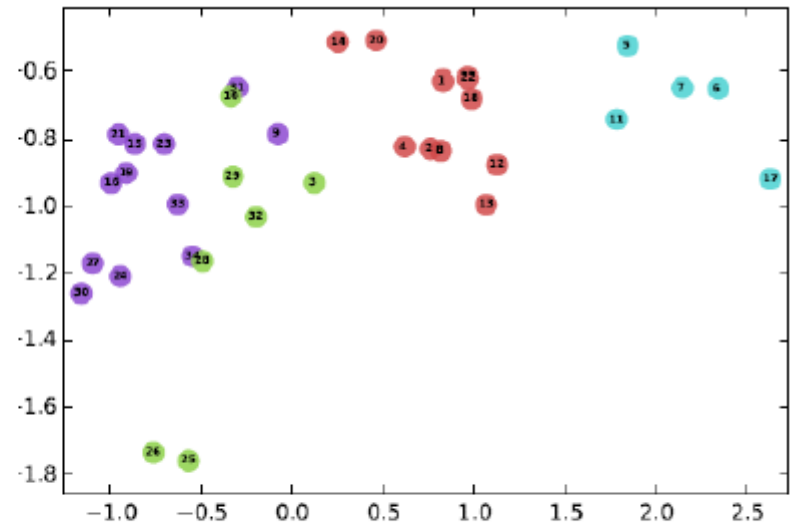


# Example

- Zachary's Karate Network:



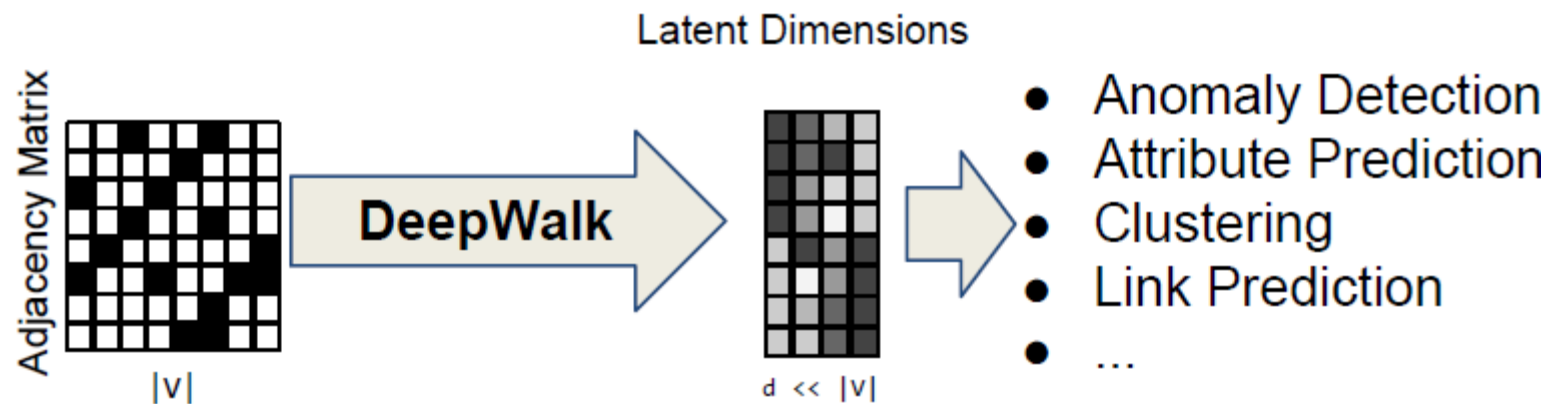
Input



Output

# DeepWalk

- DeepWalk learns a latent representation of adjacency matrices **using deep learning techniques developed for language modeling**



# Language modeling

- Learning a representation of a word from documents (word co-occurrence):
  - word2vec:  $\Phi: v \in V \mapsto \mathbb{R}^{|V| \times d}$
- The learned representations capture inherent structure
- Example:

$$\|\Phi(\textit{rose}) - \Phi(\textit{daisy})\| < \|\Phi(\textit{rose}) - \Phi(\textit{tiger})\|$$

# From language modeling to graphs

- Idea:
  - Nodes  $\leftrightarrow$  Words
  - Node sequences  $\leftrightarrow$  Sentences
- Generating node sequences:
  - Using random walks
    - short random walks = sentences
- Connection:
  - **Words frequency** in a natural language corpus follows a power law.
  - **Vertex frequency** in random walks on scale free graphs also follows a power law.

# Note 1

- How to generate random paths based on a graph
  - Weighted graph
  - Unweighted graph

$$P(c_i = x \mid c_{i-1} = v) = \begin{cases} \frac{\pi_{vx}}{Z} & \text{if } (v, x) \in E \\ 0 & \text{otherwise} \end{cases}$$

where  $\pi_{vx}$  is the unnormalized transition probability between nodes  $v$  and  $x$ , and  $Z$  is the normalizing constant.

# Note 2

- How to sample a discrete variable from a multinomial distribution

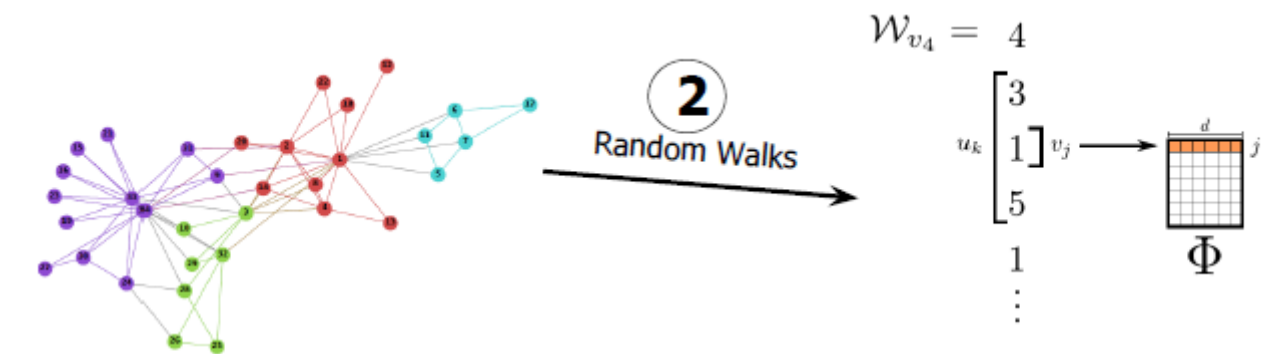


# Note 3

- What is PageRank?

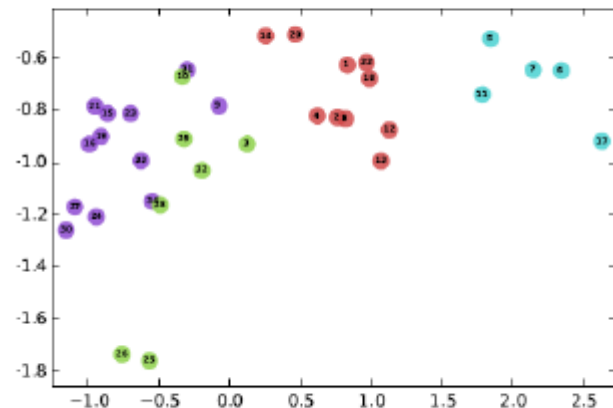
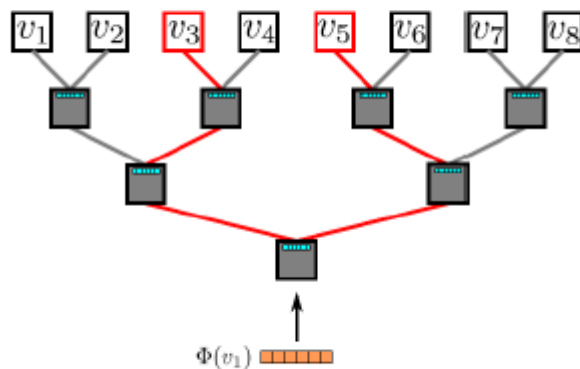
$$PR(p_i) = \frac{1-d}{N} + d \sum_{p_j \in M(p_i)} \frac{PR(p_j)}{L(p_j)}$$

# Framework



**1** Input: Graph

**3** Representation Mapping



**4** Hierarchical Softmax

**5** Output: Representation

# Representation Mapping

$$\mathcal{W}_{v_4} \equiv v_4 \rightarrow v_3 \rightarrow \textcolor{red}{v_1} \rightarrow v_5 \rightarrow v_1 \rightarrow v_{46} \rightarrow v_{51} \rightarrow v_{89}$$

$$\mathcal{W}_{v_4} = 4$$

$$u_k \begin{bmatrix} 3 \\ 1 \\ 5 \\ 1 \\ \vdots \end{bmatrix} v_j \longrightarrow \begin{array}{c} d \\ \text{Grid} \\ j \\ \Phi \end{array}$$

■ Map the vertex under focus ( $\textcolor{red}{v_1}$ ) to its representation.

■ Define a window of size  $\mathcal{W}$

■ If  $\mathcal{W} = 1$  and  $\mathcal{V} = \textcolor{red}{v_1}$

**Maximize:**  $\Pr(v_3 | \Phi(\textcolor{red}{v_1}))$

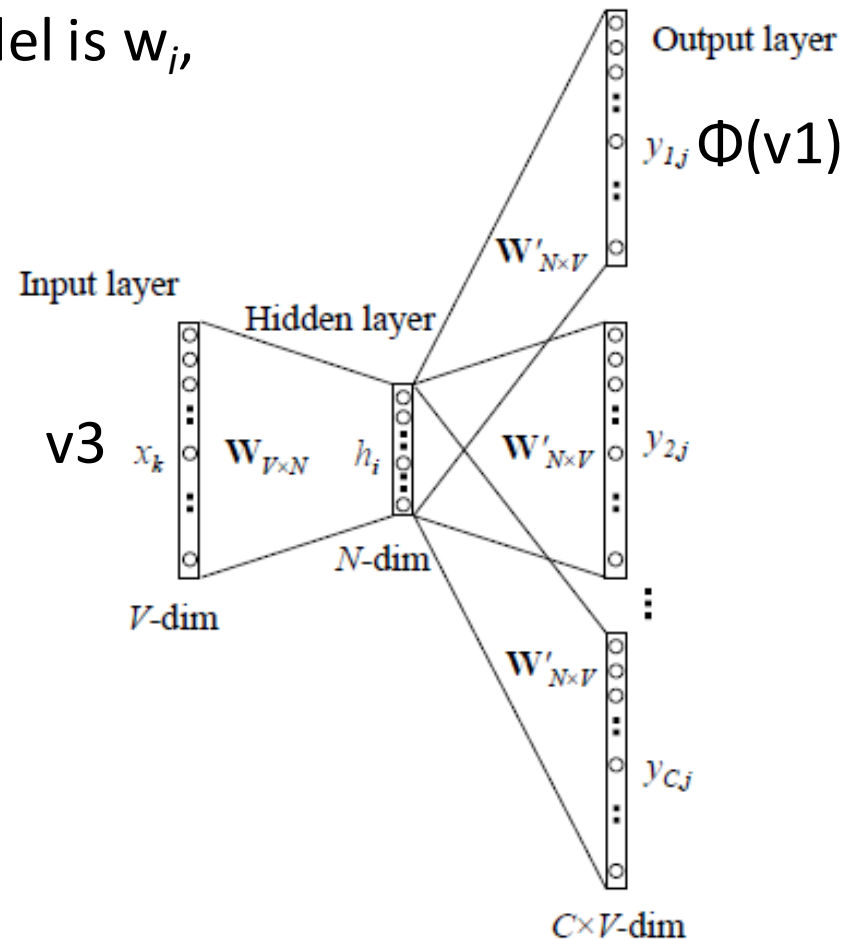
$\Pr(v_5 | \Phi(\textcolor{red}{v_1}))$

# Deep Learning Structure: Skip-gram model

Skip-gram: The input to the model is  $w_i$ ,  
and the output could be

$w_{i-1}, w_{i-2}, w_{i+1}, w_{i+2}$

**Maximize:**  $\Pr(v_3 | \Phi(v_1))$   
 $\Pr(v_5 | \Phi(v_1))$



# Experiments

- Node Classification
  - Some nodes have labels, some don't
- DataSet
  - BlogCatalog
  - Flickr
  - YouTube

# Results: BlogCatalog

	% Labeled Nodes	10%	20%	30%	40%	50%	60%	70%	80%	90%
Micro-F1(%)	DEEPWALK	<b>36.00</b>	<b>38.20</b>	<b>39.60</b>	<b>40.30</b>	<b>41.00</b>	<b>41.30</b>	41.50	41.50	42.00
	SpectralClustering	31.06	34.95	37.27	38.93	39.97	40.99	<b>41.66</b>	<b>42.42</b>	<b>42.62</b>
	EdgeCluster	27.94	30.76	31.85	32.99	34.12	35.00	34.63	35.99	36.29
	Modularity	27.35	30.74	31.77	32.97	34.09	36.13	36.08	37.23	38.18
	wvRN	19.51	24.34	25.62	28.82	30.37	31.81	32.19	33.33	34.28
	Majority	16.51	16.66	16.61	16.70	16.91	16.99	16.92	16.49	17.26
Macro-F1(%)	DEEPWALK	<b>21.30</b>	<b>23.80</b>	25.30	26.30	27.30	27.60	27.90	28.20	28.90
	SpectralClustering	19.14	23.57	<b>25.97</b>	<b>27.46</b>	<b>28.31</b>	<b>29.46</b>	<b>30.13</b>	<b>31.38</b>	<b>31.78</b>
	EdgeCluster	16.16	19.16	20.48	22.00	23.00	23.64	23.82	24.61	24.92
	Modularity	17.36	20.00	20.80	21.85	22.65	23.41	23.89	24.20	24.97
	wvRN	6.25	10.13	11.64	14.24	15.86	17.18	17.98	18.86	19.57
	Majority	2.52	2.55	2.52	2.58	2.58	2.63	2.61	2.48	2.62

# Network Embedding Models

- DeepWalk
- **Node2vec** (Grover et al., KDD 2016)
- GENE
- LINE
- SDNE

# Node2Vec

- A generalized version of DeepWalk
  - Objective function

$$\max_f \sum_{u \in V} \log \Pr(N_S(u) | f(u)).$$

- Conditional independence

$$\Pr(N_S(u) | f(u)) = \prod_{n_i \in N_S(u)} \Pr(n_i | f(u)).$$

- Symmetry in feature space

$$\Pr(n_i | f(u)) = \frac{\exp(f(n_i) \cdot f(u))}{\sum_{v \in V} \exp(f(v) \cdot f(u))}.$$



# Node2Vec

$$N_S(u) \subset V$$

- a network neighborhood of node  $u$  generated through a neighborhood sampling strategy  $S$ .
- The key lies in how to find a neighbor on the graph
- How *DeepWalk* solve this?

$$P(c_i = x \mid c_{i-1} = v) = \begin{cases} \frac{\pi_{vx}}{Z} & \text{if } (v, x) \in E \\ 0 & \text{otherwise} \end{cases}$$

where  $\pi_{vx}$  is the unnormalized transition probability between nodes  $v$  and  $x$ , and  $Z$  is the normalizing constant.

# How Node2vec Do this?

- Motivation

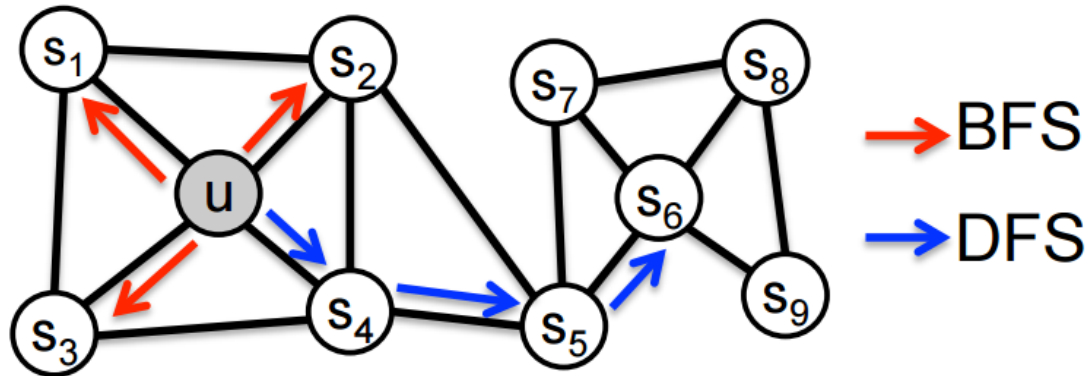


Figure 1: BFS and DFS search strategies from node  $u$  ( $k = 3$ ).

- BFS: broader → *homophily*
- DFS: deeper → *structural equivalence*

# How Node2vec Do this?

- Can we combine the merits of DFS and BFS
  - BFS: broader  $\rightarrow$  *homophily*
  - DFS: deeper  $\rightarrow$  *structural equivalence*

$$P(c_i = x \mid c_{i-1} = v) = \begin{cases} \frac{\pi_{vx}}{Z} & \text{if } (v, x) \in E \\ 0 & \text{otherwise} \end{cases}$$

$$\pi_{vx} = \alpha_{pq}(t, x) \cdot w_{vx}$$

$$\alpha_{pq}(t, x) = \begin{cases} \frac{1}{p} & \text{if } d_{tx} = 0 \\ 1 & \text{if } d_{tx} = 1 \\ \frac{1}{q} & \text{if } d_{tx} = 2 \end{cases}$$

# How Node2vec Do this?

- Explaining the sampling strategy

$$P(c_i = x \mid c_{i-1} = v) = \begin{cases} \frac{\pi_{vx}}{Z} & \text{if } (v, x) \in E \\ 0 & \text{otherwise} \end{cases}$$

$$\pi_{vx} = \alpha_{pq}(t, x) \cdot w_{vx}$$

$$\alpha_{pq}(t, x) = \begin{cases} \frac{1}{p} & \text{if } d_{tx} = 0 \\ 1 & \text{if } d_{tx} = 1 \\ \frac{1}{q} & \text{if } d_{tx} = 2 \end{cases}$$

---

**Return parameter,  $p$ .** Parameter  $p$  controls the likelihood of immediately revisiting a node in the walk.

**In-out parameter,  $q$ .** Parameter  $q$  allows the search to differentiate between “inward” and “outward” nodes.

# Node2vec Algorithm

---

**Algorithm 1** The *node2vec* algorithm.

---

**LearnFeatures** (Graph  $G = (V, E, W)$ , Dimensions  $d$ , Walks per node  $r$ , Walk length  $l$ , Context size  $k$ , Return  $p$ , In-out  $q$ )  
     $\pi = \text{PreprocessModifiedWeights}(G, p, q)$   
     $G' = (V, E, \pi)$   
    Initialize *walks* to Empty  
    **for**  $iter = 1$  **to**  $r$  **do**  
        **for all** nodes  $u \in V$  **do**  
             $walk = \text{node2vecWalk}(G', u, l)$   
            Append  $walk$  to *walks*  
     $f = \text{StochasticGradientDescent}(k, d, walks)$   
    **return**  $f$

---

**node2vecWalk** (Graph  $G' = (V, E, \pi)$ , Start node  $u$ , Length  $l$ )  
    Initialize *walk* to  $[u]$   
    **for**  $walk\_iter = 1$  **to**  $l$  **do**  
         $curr = walk[-1]$   
         $V_{curr} = \text{GetNeighbors}(curr, G')$   
         $s = \text{AliasSample}(V_{curr}, \pi)$   
        Append  $s$  to *walk*  
    **return** *walk*

---

# Comparison between DeepWalk and Node2vec

- They actually have the same objective function and formulations
- The difference lies in how to generate random walks
- BEAUTY: node  $\rightarrow$  word, path  $\rightarrow$  sentence

# Network Embedding Models

- DeepWalk
- Node2vec
- **GENE** (Chen et al., CIKM 2016)
- LINE
- SDNE

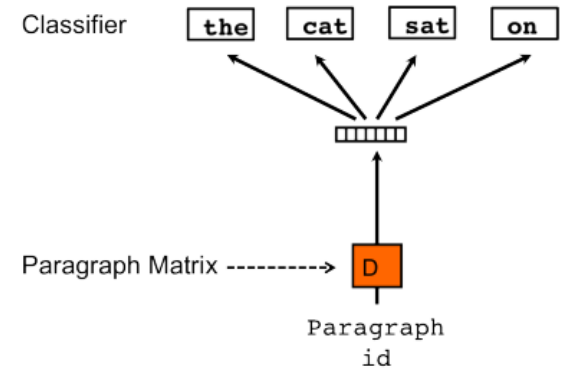
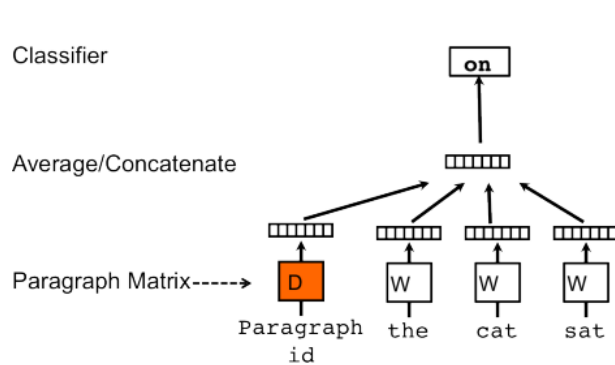
# GENE

- Incorporate Group Information to Enhance Network Embedding
  - When group information is available, how to model it?
    - Group  $\rightarrow_{\text{control}}$  member



# GENE

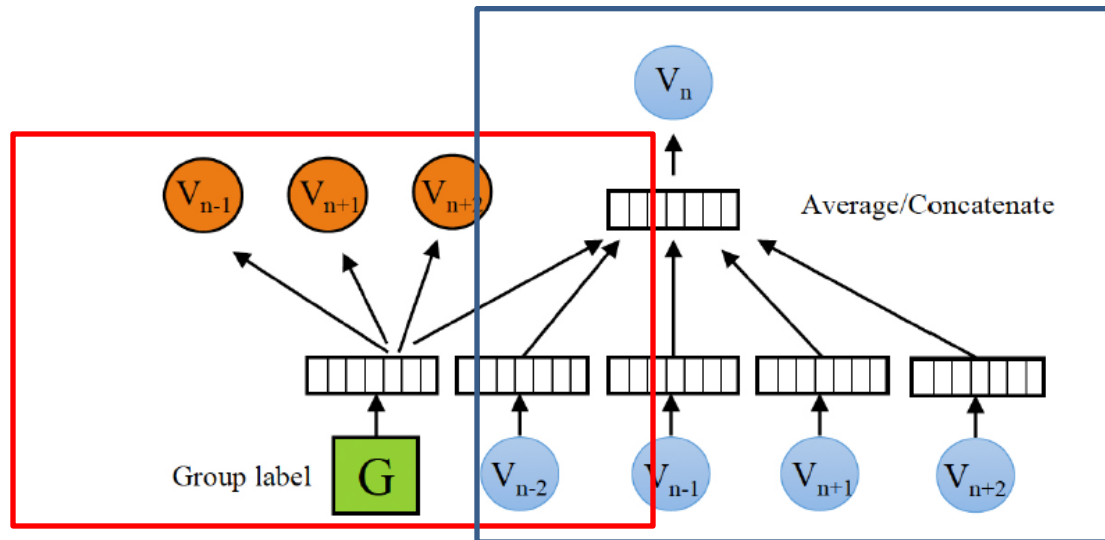
- Recall doc2vec



- How to use doc2vec to model group and member vectors

# GENE

- Incorporate Group Information to Enhance Network Embedding
  - When group information is available, how to model it?



# GENE

- Formulate the idea

$$\mathcal{L} = \sum_{g_i \in C} (\alpha \sum_{W \in W_{g_i}} \sum_{v_j \in W} \log p(v_j | v_{j-k}, \dots, v_{j+k}, g_i) + \beta \sum_{\hat{v}_j \in \hat{W}_{g_i}} \log p(\hat{v}_j | g_i)), \quad (1)$$

$$\log p(v_j | v_{j-k}, \dots, v_{j+k}, g_i) = \frac{\exp(\bar{u}^T u'_j)}{\sum_{n=1}^M \exp(\bar{u}^T u'_n)}, \quad (2)$$

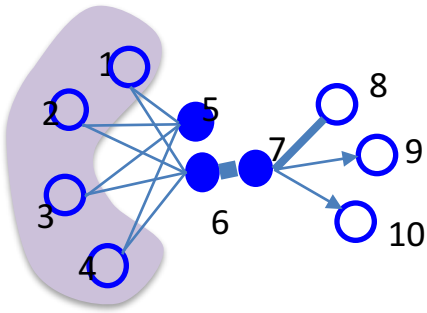
$$\log p(\hat{v}_j | g_i) = \frac{\exp(u_{g_i}^T \hat{u}_j)}{\sum_{n=1}^M \exp(u_{g_i}^T \hat{u}_n)}, \quad (3)$$

# Network Embedding Models

- DeepWalk
- Node2vec
- GENE
- **LINE** (Tang et al., WWW 2015)
- SDNE

# LINE

## First-order Proximity

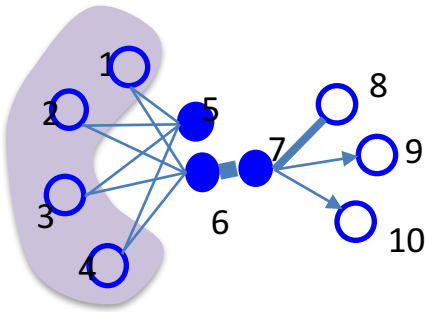


Vertex **6** and **7** have a large first-order proximity

- 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

# LINE

## Second-order Proximity



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

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

$$\hat{p}_6 = (1, 1, 1, 1, 0, 0, 5, 0, 0, 0)$$

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

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

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

# Questions

- How to characterize the first-order and second-order proximity?
  - We assume each node is associated with a low-dimensional latent factor

# LINE

## Preserving the First-order Proximity

- 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

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



# LINE

## Preserving the Second-order Proximity

- 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:

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

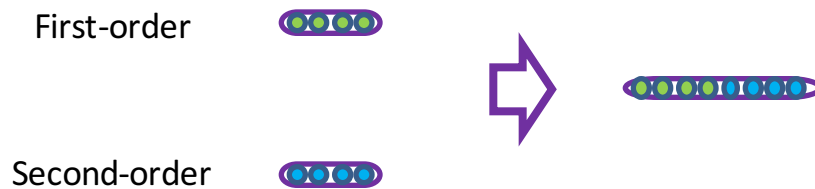
$\lambda_i$ : Prestige of vertex in the network  
 $\lambda_i = \sum_j w_{ij}$

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

# LINE

## Preserving both Proximity

- Concatenate the embeddings individually learned by the two proximity



# LINE

## Optimization

- Stochastic gradient descent + Negative Sampling
  - Randomly sample an edge and multiple negative edges
- The gradient w.r.t the embedding with edge  $(i, j)$

$$\frac{\partial O_2}{\partial \vec{u}_i} = w_{ij} \cdot \frac{\partial \log p_2(v_j | v_i)}{\partial \vec{u}_i}$$

Multiplied by the weight of the edge  $w_{ij}$

- Problematic when the weights of the edges diverge
  - The scale of the gradients with different edges diverges
- Solution: **edge sampling**
  - Sample the edges according to their weights and treat the edges as binary
- Complexity:  $O(dK|E|)$ 
  - Linear to the dimension  $d$ , the number of negative samples  $K$ , and the number of edges  $|E|$

# Questions

- How to model nodes with a small degree
- How to model new nodes

# Solutions

## Embedding Vertices of small degrees

- Sparse information in the neighborhood
- Solution: expand the neighbors by adding higher-order neighbors
  - e.g., neighbors of neighbors
  - breadth-first search
  - only consider the second-order neighbors

## Embedding New Vertices

- Fix existing embeddings, and optimize w.r.t the new ones
- Objective

$$- \sum_{j \in N(i)} w_{ji} \log p_1(v_j, v_i) \quad \text{or}$$

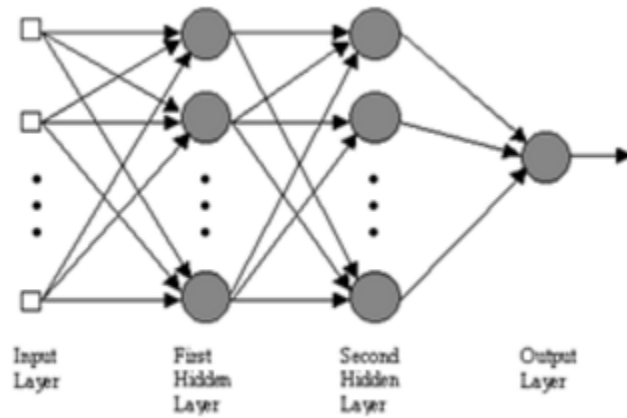
$$- \sum_{j \in N(i)} w_{ji} \log p_2(v_j | v_i)$$

# Network Embedding Models

- DeepWalk
- Node2vec
- GENE
- LINE
- **SDNE** (Wang et al., KDD 2016)

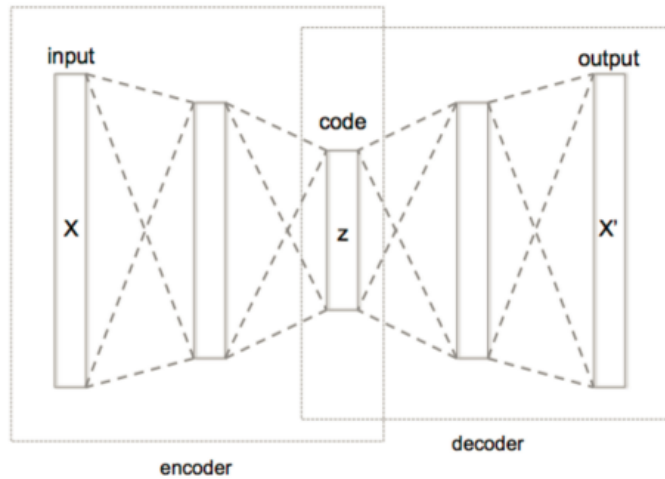
# SDNE

- Preliminary
  - Multi-layer perceptron



# SDNE

- Preliminary
  - Autoencoder



$$\phi : \mathcal{X} \rightarrow \mathcal{F}$$

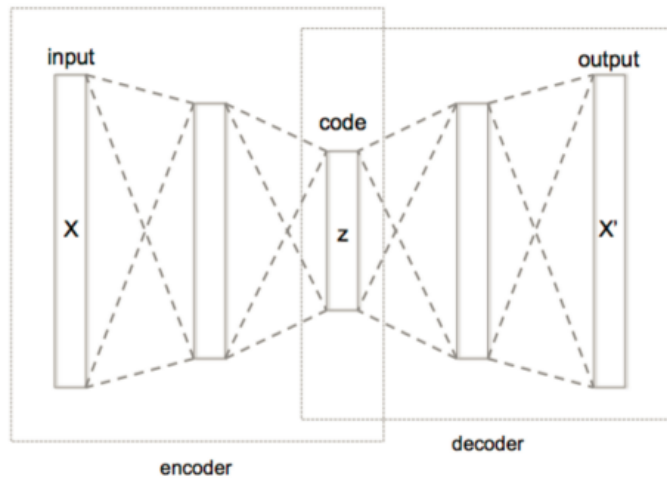
$$\psi : \mathcal{F} \rightarrow \mathcal{X}$$

$$\arg \min_{\phi, \psi} \|X - (\psi \circ \phi)X\|^2$$



# SDNE

- Preliminary
  - Autoencoder
    - The simplest case: a single hidden layer



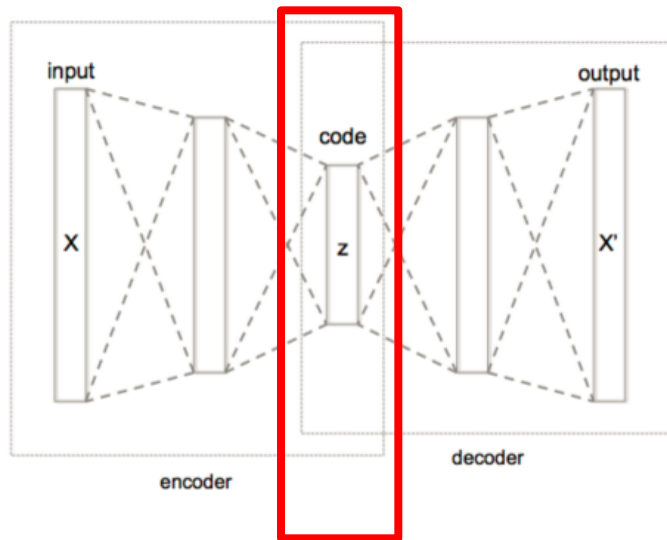
$$\mathbf{z} = \sigma_1(\mathbf{W}\mathbf{x} + \mathbf{b})$$

$$\mathbf{x}' = \sigma_2(\mathbf{W}'\mathbf{z} + \mathbf{b}')$$

$$\mathcal{L}(\mathbf{x}, \mathbf{x}') = \|\mathbf{x} - \mathbf{x}'\|^2$$

# SDNE

- Preliminary
  - Autoencoder
    - The simplest case: a single hidden layer



$$\mathbf{z} = \sigma_1(\mathbf{W}\mathbf{x} + \mathbf{b})$$

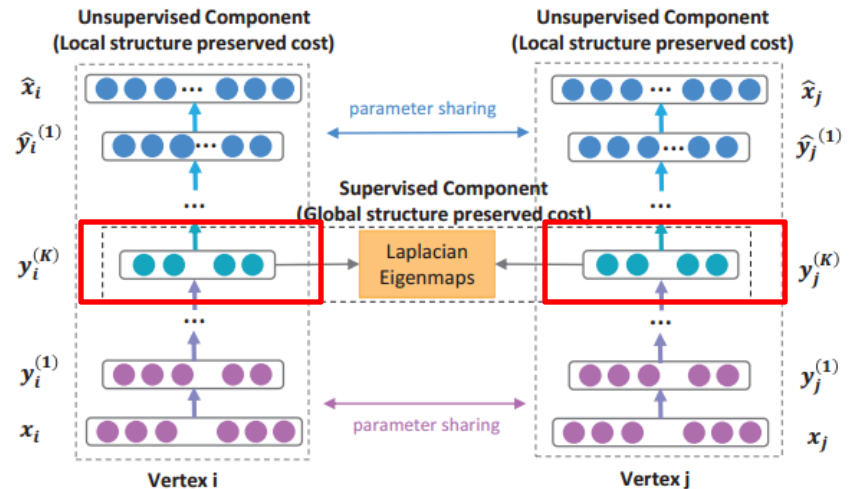
$$\mathbf{x}' = \sigma_2(\mathbf{W}'\mathbf{z} + \mathbf{b}')$$

$$\mathcal{L}(\mathbf{x}, \mathbf{x}') = \|\mathbf{x} - \mathbf{x}'\|^2$$

# SDNE

- First-order proximity
  - Linked nodes should be coded similarly

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



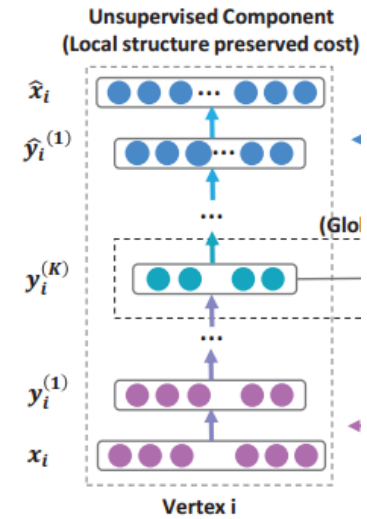
$$\mathbf{y}_i^{(1)} = \sigma(W^{(1)}\mathbf{x}_i + \mathbf{b}^{(1)})$$

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

# SDNE

- Second-order proximity
  - The model should reconstruct the neighborhood vectors
  - Similar nodes even without links can have similar codes
    - Or we can not reconstruct the neighborhood

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



$$\begin{aligned}\mathbf{y}_i^{(1)} &= \sigma(W^{(1)}\mathbf{x}_i + \mathbf{b}^{(1)}) \\ \mathbf{y}_i^{(k)} &= \sigma(W^{(k)}\mathbf{y}_i^{(k-1)} + \mathbf{b}^{(k)}), k = 2, \dots, K\end{aligned}$$

# SDNE

- Network reconstruction

**Table 4: MAP on ARXIV-GRQC and BLOGCATALOG on reconstruction task**

Method	ARXIV-GRQC					BLOGCATALOG				
	<i>SDNE</i>	<i>GraRep</i>	<i>LINE</i>	<i>DeepWalk</i>	<i>LE</i>	<i>SDNE</i>	<i>GraRep</i>	<i>LINE</i>	<i>DeepWalk</i>	<i>LE</i>
MAP	<b>0.836**</b>	0.05	0.69	0.58	0.23	<b>0.63**</b>	0.42	0.58	0.28	0.12

Significantly outperforms GraRep at the: \*\* 0.01 level.

- Link prediction

**Table 5: *precision@k* on ARXIV GR-QC for link prediction**

Algorithm	<i>P@2</i>	<i>P@10</i>	<i>P@100</i>	<i>P@200</i>	<i>P@300</i>	<i>P@500</i>	<i>P@800</i>	<i>P@1000</i>	<i>P@10000</i>
<i>SDNE</i>	<b>1</b>	<b>1</b>	<b>1</b>	<b>1</b>	<b>1*</b>	<b>0.99**</b>	<b>0.97**</b>	<b>0.91**</b>	<b>0.257**</b>
<i>LINE</i>	1	1	1	1	0.99	0.936	0.74	0.79	0.2196
<i>DeepWalk</i>	1	0.8	0.6	0.555	0.443	0.346	0.2988	0.293	0.1591
<i>GraRep</i>	1	0.2	0.04	0.035	0.033	0.038	0.035	0.035	0.019
<i>Common Neighbor</i>	1	1	1	0.96	0.9667	0.98	0.8775	0.798	0.192
<i>LE</i>	1	1	0.93	0.855	0.827	0.66	0.468	0.391	0.05

Significantly outperforms Line at the: \*\* 0.01 and \* 0.05 level, paired t-test.

# Network Embedding Models

- DeepWalk
  - Node sentences + word2vec
- Node2vec
  - DeepWalk + more sampling strategies
- GENE
  - Group~document + doc2vec(DM, DBOW)
- LINE
  - Shallow + first-order + second-order proximity
- SDNE
  - Deep + First-order + second-order proximity

# Applications of Network Embedding

- Basic applications
- Data Visualization
- Text classification
- Recommendation

# Basic Applications

- Network reconstruction
- Link prediction
- Clustering
- Feature coding
  - Node classification
    - Demographic prediction
    - **Question: how to infer the demographics of a microblog user**



# Applications of Network Embedding

- Basic applications
- **Data Visualization** (Tang et al., WWW 2016)
- Text classification
- Recommendation

# Data Visualization

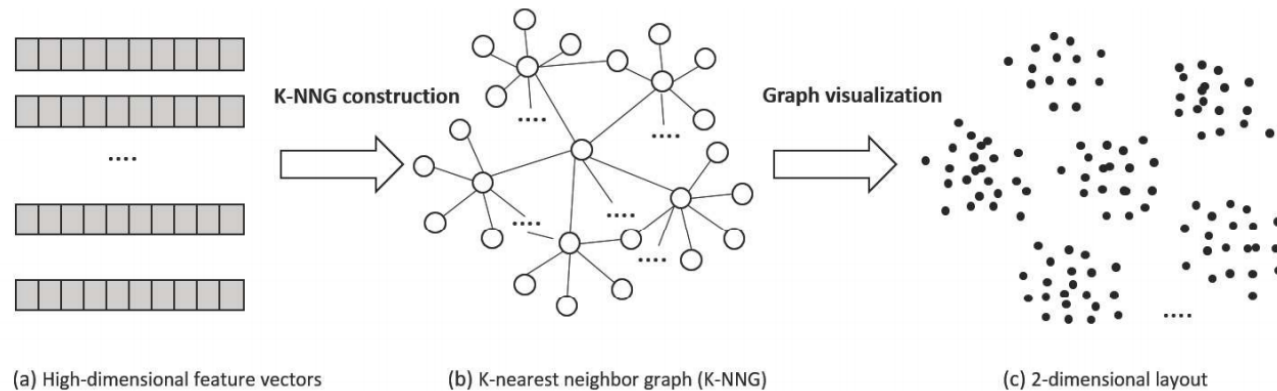


Figure 1: A typical pipeline of data visualization by first constructing a K-nearest neighbor graph and then projecting the graph into a low-dimensional space.

# Data Visualization

- Construction of the KNN graph

For the weights of the edges in the K-nearest neighbor graph, we use the same approach as t-SNE. The conditional probability from data  $\vec{x}_i$  to  $\vec{x}_j$  is first calculated as:

$$p_{j|i} = \frac{\exp(-\|\vec{x}_i - \vec{x}_j\|^2 / 2\sigma_i^2)}{\sum_{(i,k) \in E} \exp(-\|\vec{x}_i - \vec{x}_k\|^2 / 2\sigma_i^2)}, \text{ and} \quad (1)$$
$$p_{i|i} = 0,$$

Then the graph is symmetrized through setting the weight between  $\vec{x}_i$  and  $\vec{x}_j$  as:

$$w_{ij} = \frac{p_{j|i} + p_{i|j}}{2N}. \quad (2)$$

# Data Visualization

- Visualization-based embedding

$$P(e_{ij} = 1) = f(\|\vec{y}_i - \vec{y}_j\|),$$

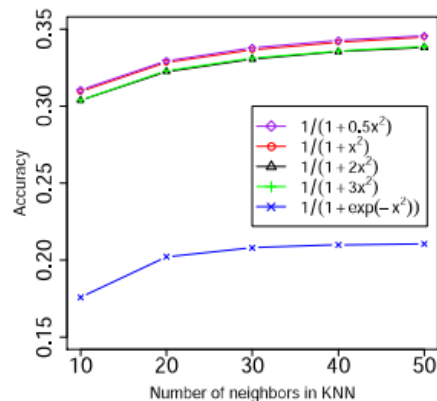
$$P(e_{ij} = w_{ij}) = P(e_{ij} = 1)^{w_{ij}}.$$

$$\begin{aligned} O &= \prod_{(i,j) \in E} p(e_{ij} = 1)^{w_{ij}} \prod_{(i,j) \in \bar{E}} (1 - p(e_{ij} = 1))^\gamma \\ &\propto \sum_{(i,j) \in E} w_{ij} \log p(e_{ij} = 1) + \sum_{(i,j) \in \bar{E}} \gamma \log(1 - p(e_{ij} = 1)), \end{aligned}$$

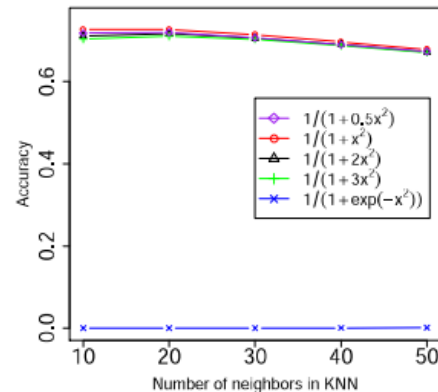
# Data Visualization

- Non-linear function

$$P(e_{ij} = 1) = f(\|\vec{y}_i - \vec{y}_j\|),$$



(a) WikiDoc



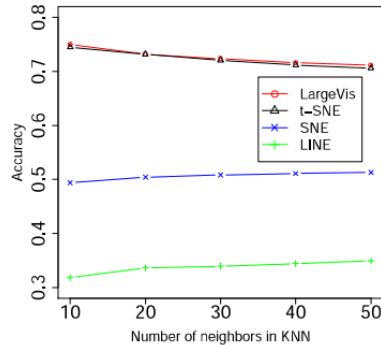
(b) LiveJournal

Figure 4: Comparing different probabilistic functions.

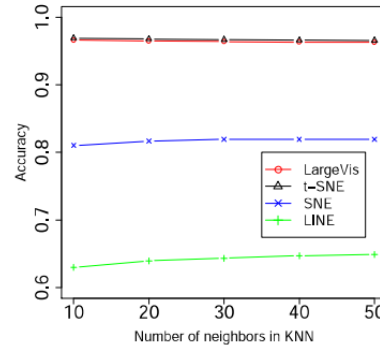
$$f(x) = \frac{1}{1+x^2}$$

# Data Visualization

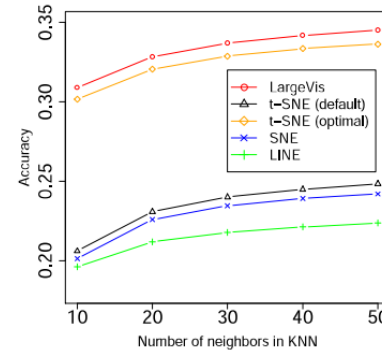
- Accuracy



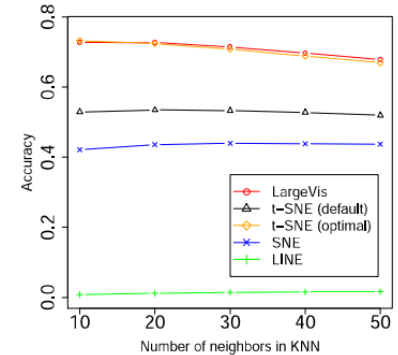
(a) 20NG



(b) MNIST



(c) WikiDoc



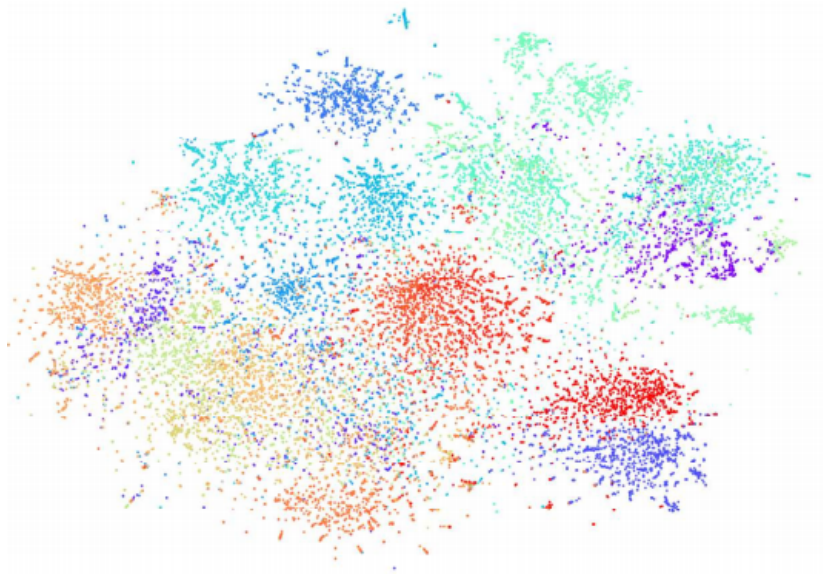
(d) LiveJournal

- Running time

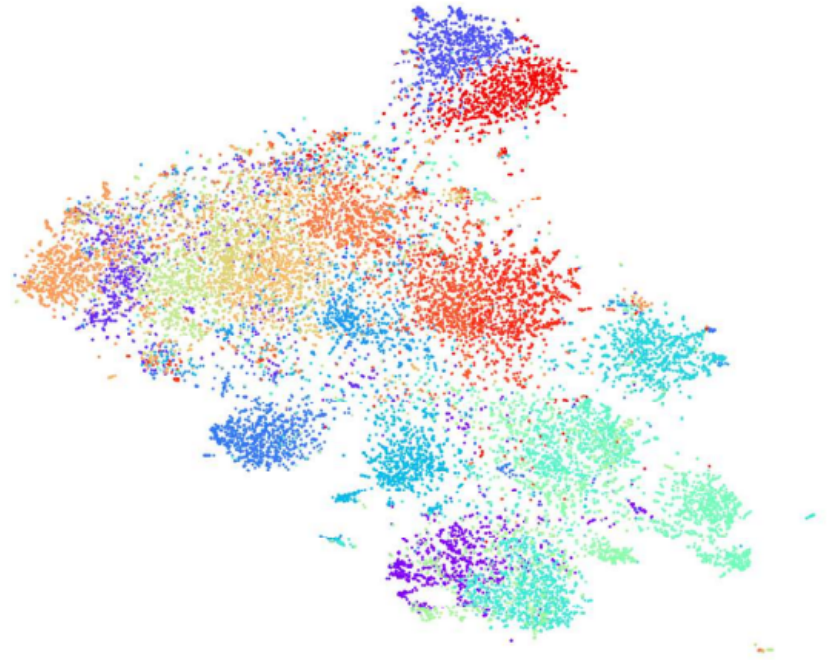
Table 2: Comparison of running time (hours) in graph visualization between the t-SNE and LargeVis.

Algorithm	20NG	MNIST	WikiWord	WikiDoc	LiveJournal	CSAuthor	DBLPPaper
t-SNE	0.12	0.41	9.82	45.01	70.35	28.33	18.73
LargeVis	0.14	0.23	2.01	5.60	9.26	4.24	3.19
Speedup Rate	0	0.7	3.9	7	6.6	5.7	4.9

# Data Visualization



(a) 20NG (t-SNE)



(b) 20NG (LargeVis)

# Applications of Network Embedding

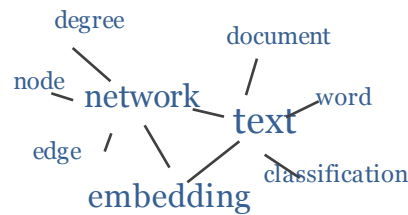
- Basic applications
- Data Visualization
- **Text classification** (Tang et al., KDD 2015)
- Recommendation



# Text Classification

## Network embedding helps text modeling

Text representation, e.g., word and document representation, ...  
Deep learning has been attracting increasing attention ...  
A future direction of deep learning is to integrate unlabeled data ...  
...  
The Skip-gram model is quite effective and efficient ...  
Information networks encode the relationships between the data objects ...



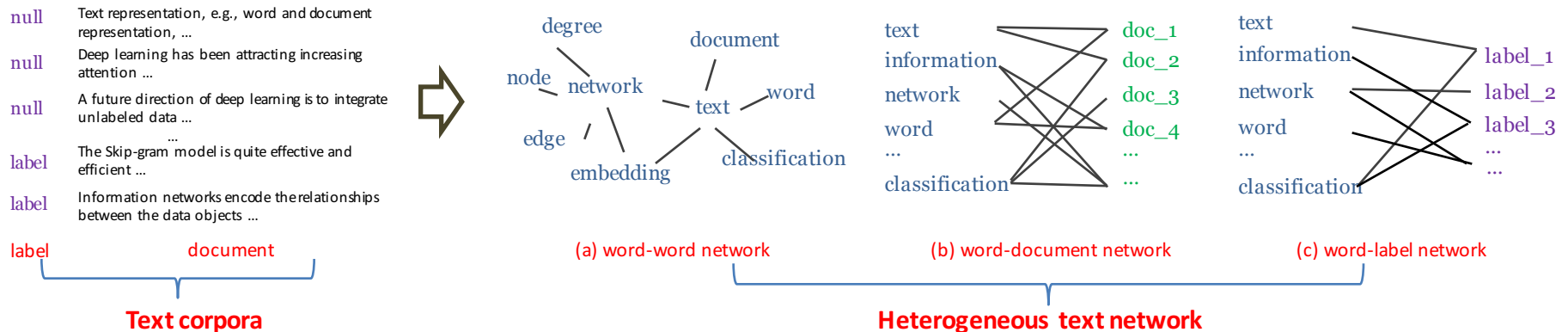
**Free text**

**word co-occurrence network**

If we have the word network, we can a network embedding model to learn word representations.

# Text Classification

- Adapt the advantages of unsupervised text embedding approaches but naturally utilize the *labeled* data for specific tasks
- Different levels of word co-occurrences: *local context-level, document-level, label-level*



# Text Classification

## Bipartite Network Embedding

- Extend previous work **LINE** (Tang et al. WWW'2015) on large-scale information network embedding
  - Preserve the *first-order* and *second-order* proximity
  - Only consider the *second-order* proximity here

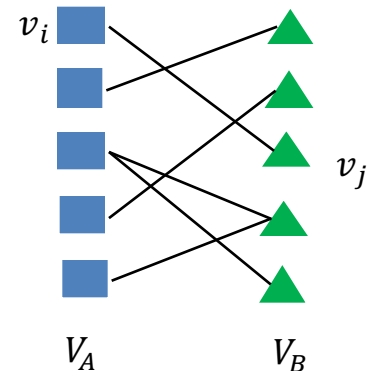
- For each edge  $(v_i, v_j)$ , define a conditional probability

$$p(v_j|v_i) = \frac{\exp(\vec{u}_j^T \cdot \vec{u}_i)}{\sum_{j' \in B} \exp(\vec{u}_{j'}^T \cdot \vec{u}_i)}$$

- Objective:

$$O = - \sum_{(i,j) \in E} w_{ij} \log p(v_j|v_i)$$

- Edge sampling and negative sampling for optimization



Tang et al. **LINE: Large-scale Information Network Embedding**. WWW'2015

# Text Classification

## Heterogeneous Text Network Embedding

- Heterogeneous text network: three bipartite networks
  - Word-word (word-context), word-document, word-label network
  - Jointly embed the three bipartite networks
- Objective

$$O_{pte} = O_{ww} + O_{wd} + O_{wl}$$

- where

$$O_{ww} = - \sum_{(i,j) \in E_{ww}} w_{ij} \log p(v_i | v_j)$$

Objective for **word-word** network

$$O_{wd} = - \sum_{(i,j) \in E_{wd}} w_{ij} \log p(v_i | d_j)$$

Objective for **word-document** network

$$O_{wl} = - \sum_{(i,j) \in E_{wl}} w_{ij} \log p(v_i | l_j)$$

Objective for **word-label** network

# Text Classification

## Results on Long Documents: Predictive

		20newsgroup		Wikipedia		IMDB	
Type	Algorithm	Micro-F1	Macro-F1	Micro-F1	Macro-F1	Micro-F1	Macro-F1
Unsupervised	LINE( $G_{wd}$ )	79.73	78.40	80.14	80.13	89.14	89.14
Predictive embedding	CNN	78.85	78.29	79.72	79.77	86.15	86.15
	CNN(pretrain)	80.15	79.43	79.25	79.32	89.00	89.00
	PTE( $G_{wl}$ )	82.70	81.97	79.00	79.02	85.98	85.98
	PTE( $G_{ww} + G_{wl}$ )	83.90	83.11	81.65	81.62	89.14	89.14
	PTE( $G_{wd} + G_{wl}$ )	84.39	83.64	82.29	82.27	89.76	89.76
	PTE(pretrain)	82.86	82.12	79.18	79.21	86.28	86.28
	PTE(joint)	84.20	83.39	82.51	82.49	89.80	89.80

PTE(joint) > PTE(pretrain)

PTE(joint) > PTE( $G_{wl}$ )

PTE(joint) > CNN/CNN(pretrain)

# Text Classification

## Results on **Short** Documents: Predictive

		DBLP		MR		Twitter	
Type	Algorithm	Micro-F1	Macro-F1	Micro-F1	Macro-F1	Micro-F1	Macro-F1
Unsupervised embedding	LINE ( $G_{ww} + G_{wd}$ )	74.22	70.12	71.13	71.12	73.84	73.84
Predictive embedding	CNN	76.16	73.08	72.71	72.69	<b>75.97</b>	<b>75.96</b>
	CNN(pretrain)	75.39	72.28	68.96	68.87	75.92	75.92
	PTE( $G_{wl}$ )	76.45	72.74	73.44	73.42	73.92	73.91
	PTE( $G_{ww} + G_{wl}$ )	76.80	73.28	72.93	72.92	74.93	74.92
	PTE( $G_{wd} + G_{wl}$ )	<b>77.46</b>	<b>74.03</b>	73.13	73.11	75.61	75.61
	PTE(pretrain)	76.53	72.94	73.27	73.24	73.79	73.79
	PTE(joint)	77.15	73.61	<b>73.58</b>	<b>73.57</b>	75.21	75.21

PTE(joint) > PTE(pretrain)

PTE(joint) > PTE( $G_{wl}$ )

PTE(joint)  $\approx$  CNN/CNN(pretrain)

# Applications of Network Embedding

- Basic applications
- Data Visualization
- Text classification
- **Recommendation** (Zhao et al., AIRS 2016)

# Traditional Recommendation Methods

- Popularity
- Item-KNN
- User-KNN
- Mixture of item-KNN and user-KNN
- Bayesian Personalized Ranking



# Recommendation

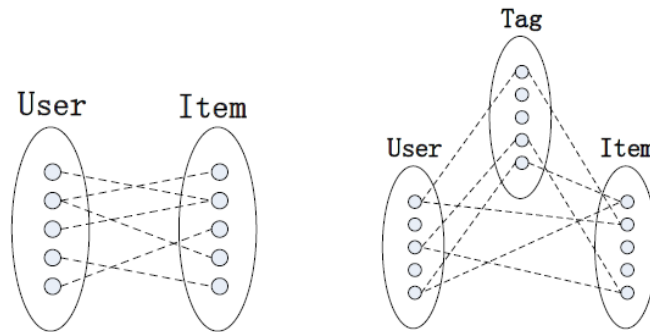
- From training records to networks

**Definition 3. Bipartite User-Item (UI) Network.** Let  $\mathcal{U}$  denote the set of all the users, and  $\mathcal{I}$  denote the set of all the items. A bipartite user-item network can be denoted by  $\mathcal{G}^{(bi)} = (\mathcal{V}, \mathcal{E}, \mathbf{W})$ , where the vertex set  $\mathcal{V} = \mathcal{U} \cup \mathcal{I}$ , the edge set  $\mathcal{E} \subset \mathcal{U} \times \mathcal{I}$ , the weight matrix  $\mathbf{W}$  stores the edge weights, and  $W_{u,i}$  denote the link weight between a user  $u$  and an item  $i$ .

**Definition 4. Tripartite User-Item-Tag (UIT) Network.** Let  $\mathcal{U}$  denote the set of all the users,  $\mathcal{I}$  denote the set of all the items, and  $\mathcal{T}$  denote the set of all the tags. A tripartite user-item-tag network can be denoted by  $\mathcal{G}^{(tri)} = (\mathcal{V}, \mathcal{E}, \mathbf{W})$ , where the vertex set  $\mathcal{V} = \mathcal{U} \cup \mathcal{I} \cup \mathcal{T}$ , the edge set  $\mathcal{E} \subset ((\mathcal{U} \times \mathcal{I}) \cup (\mathcal{U} \times \mathcal{T}) \cup (\mathcal{I} \times \mathcal{T}))$ , and the weight matrix  $\mathbf{W}$  stores the edge weights.

# Recommendation

- Learning Distributed Representations for Recommender Systems with a Network Embedding Approach
  - Motivation



(a) User-item bipartite network. (b) User-item-tag tripartite network.

# Recommendation

- Given any edge in the network

$$P(e_s, e_t) = \sigma(\mathbf{v}_{e_s}^\top \cdot \mathbf{v}_{e_t}) = \frac{1}{1 + \exp(-\mathbf{v}_{e_s}^\top \cdot \mathbf{v}_{e_t})}.$$

$$\hat{P}(e_s, e_t) = \frac{W_{e_s, e_t}}{\sum_{(e_{s'}, e_{t'}) \in \mathcal{E}} W_{e_{s'}, e_{t'}}}.$$

$$L(\mathcal{G}) = D_{KL}(\hat{P}(\cdot, \cdot) || P(\cdot, \cdot)) \propto \sum_{(e_s, e_t) \in \mathcal{E}} W_{e_s, e_t} \log P(e_s, e_t).$$

# Recommendation

- User-item recommendation

**Table 2.** Performance comparisons of the proposed method and baselines on item recommendation.

Methods	JD				MovieLens			
	P@10	R@10	MAP	MRR	P@10	R@10	MAP	MRR
BPR	0.171	0.360	0.337	0.564	0.097	0.169	0.148	0.195
DeepWalk	0.259	0.443	0.502	0.806	0.203	0.243	0.249	0.358
NERM	<b>0.275</b>	<b>0.477</b>	<b>0.528</b>	<b>0.819</b>	<b>0.206</b>	<b>0.256</b>	<b>0.258</b>	<b>0.368</b>

# Recommendation

- User-item-tag recommendation

**Table 4.** Performance comparisons of the proposed methods and baselines on tag recommendation.

Methods	Last.fm						Bookmarks					
	P@1	R@1	F@1	P@5	R@5	F@5	P@1	R@1	F@1	P@5	R@5	F@5
PITF	0.305	0.125	0.178	0.189	0.351	0.245	0.381	0.132	0.197	0.204	0.304	0.244
DeepWalk	0.088	0.044	0.059	0.040	0.099	0.057	0.064	0.024	0.035	0.038	0.074	0.050
NERM	<b>0.327</b>	<b>0.165</b>	<b>0.220</b>	0.182	<b>0.370</b>	0.244	<b>0.396</b>	<b>0.135</b>	<b>0.201</b>	<b>0.228</b>	<b>0.323</b>	<b>0.267</b>

# Conclusions

- There are no boundaries between data types and research areas in terms of mythologies
  - Data models are the core
- Even if the ideas are similar, we can move from shallow to deep if the performance actually improves

# Disclaimer

- For convenience, I directly copy some original slides or figures from the referred papers. I am sorry but I did not ask for the permission of each referred author. I thank you for these slides. I will not distribute your original slides.

# References

- Tomas Mikolov, Ilya Sutskever, Kai Chen, Gregory S. Corrado, Jeffrey Dean. Distributed Representations of Words and Phrases and their Compositionality. NIPS 2013: 3111-3119
- Bryan Perozzi, Rami Al-Rfou', Steven Skiena. DeepWalk: online learning of social representations. KDD 2014: 701-710
- Aditya Grover, Jure Leskovec. node2vec: Scalable Feature Learning for Networks. KDD 2016: 855-864
- Jifan Chen, Qi Zhang, Xuanjing Huang, Incorporate Group Information to Enhance Network Embedding. CIKM 2016
- Jian Tang, Meng Qu, Mingzhe Wang, Ming Zhang, Jun Yan, Qiaozhu Mei. LINE: Large-scale Information Network Embedding. WWW 2015: 1067-1077
- Daixin Wang, Peng Cui, Wenwu Zhu. Structural Deep Network Embedding. KDD 2016: 1225-1234
- Jian Tang, Jingzhou Liu, Ming Zhang, Qiaozhu Mei. Visualizing Large-scale and High-dimensional Data. WWW 2016: 287-297
- Jian Tang, Meng Qu, Qiaozhu Mei. PTE: Predictive Text Embedding through Large-scale Heterogeneous Text Networks. KDD 2015: 1165-1174
- Wayne Xin Zhao, Jin Huang, Ji-Rong Wen. Learning Distributed Representations for Recommender Systems with a Network Embedding Approach. AIRS 2016.



# Advanced Readings

- Shiyu Chang, Wei Han, Jiliang Tang, Guo-Jun Qi, Charu C. Aggarwal, Thomas S. Huang: Heterogeneous Network Embedding via Deep Architectures. KDD 2015: 119-128
- Mingdong Ou, Peng Cui, Jian Pei, Ziwei Zhang, Wenwu Zhu: Asymmetric Transitivity Preserving Graph Embedding. KDD 2016: 1105-1114
- Thomas N. Kipf, Max Welling: Semi-Supervised Classification with Graph Convolutional Networks. CoRR abs/1609.02907 (2016)
- Mikael Henaff, Joan Bruna, Yann LeCun: Deep Convolutional Networks on Graph-Structured Data. CoRR abs/1506.05163 (2015)