

NetSMF: Large-Scale Network Embedding as Sparse Matrix Factorization

Jiezhong Qiu

Tsinghua University

June 17, 2019

Joint work with Yuxiao Dong (MSR), Hao Ma (Facebook AI),
Jian Li (IIIS, Tsinghua), Chi Wang (MSR), Kuansan Wang (MSR),
and Jie Tang (DCST, Tsinghua)

Motivation and Problem Formulation

Problem Formulation

Give a network $G = (V, E)$, aim to learn a function $f : V \rightarrow \mathbb{R}^p$ to capture neighborhood similarity and community membership.

Applications:

- ▶ link prediction
- ▶ community detection
- ▶ label classification

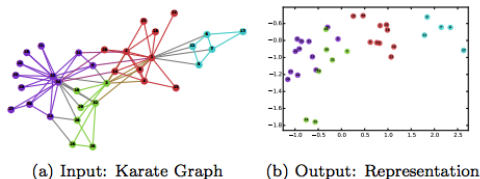


Figure 1: A toy example (Figure from DeepWalk).

Two Genres of Network Embedding Algorithm

- ▶ Local Context Methods:
 - ▶ LINE, DeepWalk, node2vec, metapath2vec.
 - ▶ Usually be formulated as a skip-gram-like problem, and optimized by SGD.
- ▶ Global Matrix Factorization Methods.
 - ▶ NetMF, GraRep, HOPE.
 - ▶ Leverage global statistics of the input networks.
 - ▶ Not necessarily a gradient-based optimization problem.
 - ▶ Usually requires explicit construction of the matrix to be factorized.

Notations

Consider an undirected weighted graph $G = (V, E)$, where $|V| = n$ and $|E| = m$.

- ▶ Adjacency matrix $\mathbf{A} \in \mathbb{R}_+^{n \times n}$:

$$\mathbf{A}_{i,j} = \begin{cases} a_{i,j} > 0 & (i,j) \in E \\ 0 & (i,j) \notin E \end{cases}.$$

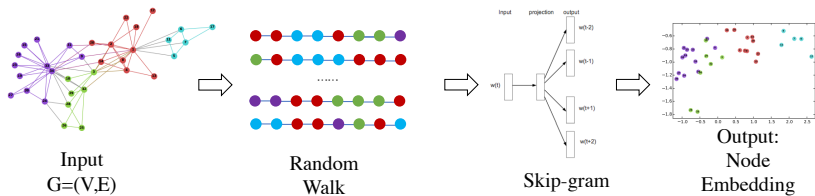
- ▶ Degree matrix $\mathbf{D} = \text{diag}(d_1, \dots, d_n)$, where d_i is the generalized degree of vertex i .
- ▶ Volume of the graph G : $\text{vol}(G) = \sum_i \sum_j \mathbf{A}_{i,j}$.

Revisit DeepWalk and NetMF

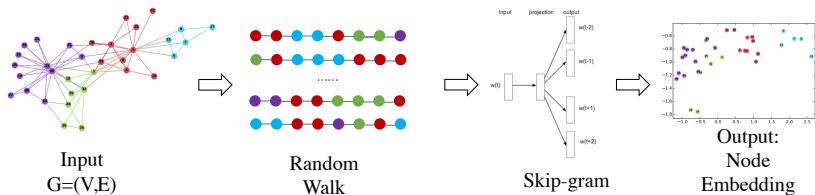
NetSMF: Network Embedding as Sparse Matrix Factorization

Experimental Results

DeepWalk and NetMF



DeepWalk and NetMF



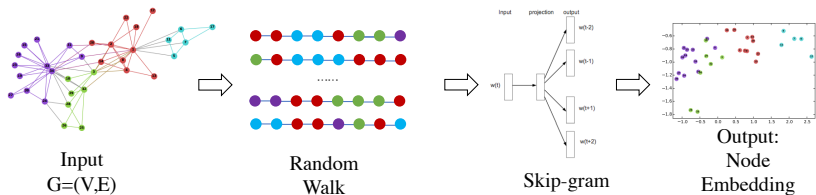
Levy & Goldberg (NIPS 14)

$$\log \left(\frac{\#(w, c) |\mathcal{D}|}{b \#(w) \#(c)} \right)$$

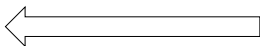
b Number of negative samples

$\#(w, c)$ Co-occurrence of w and c $\#(w)$ Occurrence of word w
 $|\mathcal{D}|$ Total number of word-context pairs $\#(c)$ Occurrence of context c

DeepWalk and NetMF



Levy & Goldberg (NIPS 14)



$$\log \left(\frac{\#(w, c) |\mathcal{D}|}{b \#(w) \#(c)} \right)$$

b Number of negative samples

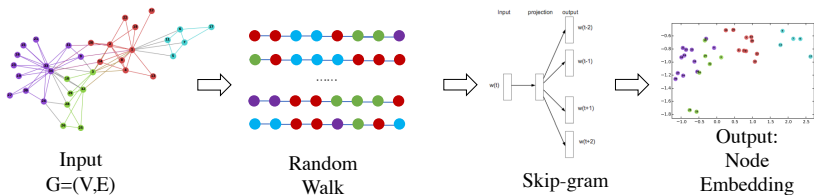
$\#(w, c)$ Co-occurrence of w and c

$\#(w)$ Occurrence of word w

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

$\#(c)$ Occurrence of context c

DeepWalk and NetMF

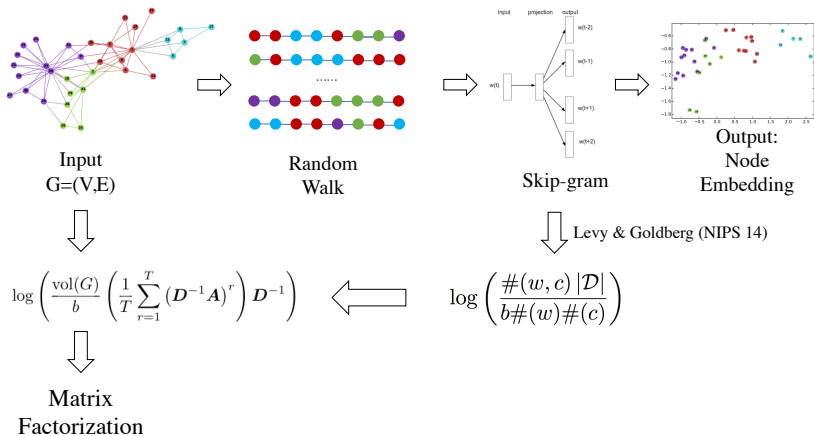


$$\log \left(\frac{\text{vol}(G)}{b} \left(\frac{1}{T} \sum_{r=1}^T (D^{-1}A)^r \right) D^{-1} \right)$$

$$\log \left(\frac{\#(w, c) |\mathcal{D}|}{b \#(w) \#(c)} \right)$$

- A** Adjacency matrix $\text{vol}(G) = \sum_i \sum_j A_{i,j}$
D Degree matrix b Number of negative samples

DeepWalk and NetMF



Revisit DeepWalk and NetMF

NetSMF: Network Embedding as Sparse Matrix Factorization

Experimental Results

Contents

Revisit DeepWalk and NetMF

NetSMF: Network Embedding as Sparse Matrix Factorization

Experimental Results

Computation Challenges of NetMF

For small world networks,

$$\frac{\text{vol}(G)}{b} \left(\frac{1}{T} \sum_{r=1}^T \underbrace{(D^{-1}A)^r}_{\text{matrix power}} \right) D^{-1} \text{ is always a dense matrix .}$$

Computation Challenges of NetMF

For small world networks,

$$\frac{\text{vol}(G)}{b} \left(\frac{1}{T} \sum_{r=1}^T \underbrace{(D^{-1}A)^r}_{\text{matrix power}} \right) D^{-1} \text{ is always a dense matrix .}$$

Why?

- ▶ In small world networks, each pair of vertices (i, j) can reach each other in a small number of hops.
- ▶ Make the corresponding matrix entry a positive value.

Computation Challenges of NetMF

For small world networks,

$$\frac{\text{vol}(G)}{b} \left(\frac{1}{T} \sum_{r=1}^T \underbrace{(D^{-1}A)^r}_{\text{matrix power}} \right) D^{-1} \text{ is always a dense matrix .}$$

Why?

- ▶ In small world networks, each pair of vertices (i, j) can reach each other in a small number of hops.
- ▶ Make the corresponding matrix entry a positive value.

Idea

- ▶ Sparse matrix is easier to handle.
- ▶ Can we achieve a matrix sparse but 'good enough' matrix.

Observation

Definition

For $\sum_{r=1}^T \alpha_r = 1$ and α_r non-negative,

$$\mathbf{L} = \mathbf{D} - \sum_{r=1}^T \alpha_r \mathbf{D} (\mathbf{D}^{-1} \mathbf{A})^r \quad (1)$$

is a T -degree random-walk matrix polynomial.

Observation

For $\alpha_1 = \dots = \alpha_T = \frac{1}{T}$:

$$\begin{aligned} & \log^\circ \left(\frac{\text{vol}(G)}{b} \left(\frac{1}{T} \sum_{r=1}^T (\mathbf{D}^{-1} \mathbf{A})^r \right) \mathbf{D}^{-1} \right) \\ &= \log^\circ \left(\frac{\text{vol}(G)}{b} \mathbf{D}^{-1} (\mathbf{D} - \mathbf{L}) \mathbf{D}^{-1} \right) \\ &\approx \log^\circ \left(\frac{\text{vol}(G)}{b} \mathbf{D}^{-1} (\mathbf{D} - \tilde{\mathbf{L}}) \mathbf{D}^{-1} \right) \end{aligned}$$

Random-walk Matrix Polynomial Sparsification

Theorem

[CCL⁺15] For random-walk matrix polynomial $\mathbf{L} = \mathbf{D} - \sum_{r=1}^T \alpha_r \mathbf{D} (\mathbf{D}^{-1} \mathbf{A})^r$, one can construct, in time $O(T^2 m \epsilon^{-2} \log^2 n)$, a $(1 + \epsilon)$ -spectral sparsifier, $\tilde{\mathbf{L}}$, with $O(n \log n \epsilon^{-2})$ non-zeros. For unweighted graphs, the time complexity can be reduced to $O(T^2 m \epsilon^{-2} \log n)$.

NetSMF—Algorithm

The proposed NetSMF algorithm consists of three steps:

- ▶ Construct a random walk matrix polynomial sparsifier, $\tilde{\mathbf{L}}$, by calling PathSampling algorithm proposed in [CCL⁺15].
- ▶ Construct a NetMF matrix sparsifier.

$$\text{trunc_log}^\circ \left(\frac{\text{vol}(G)}{b} \mathbf{D}^{-1} (\mathbf{D} - \tilde{\mathbf{L}}) \mathbf{D}^{-1} \right)$$

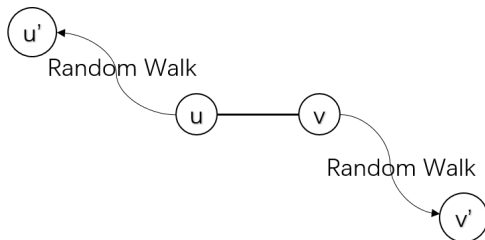
- ▶ Truncated randomized singular value decomposition.

Detailed Algorithm

Algorithm Details

PathSampling:

- ▶ Sample an edge (u, v) from edge set.
- ▶ Start very short random walk from u and arrive u' .
- ▶ Start very short random walk from v and arrive v' .
- ▶ Record vertex pair (u', v') .



Randomized SVD:

- ▶ Project origin matrix to low dimensional space by Gaussian random matrix.
- ▶ Deal with the projected small matrix.

NetSMF — System Design

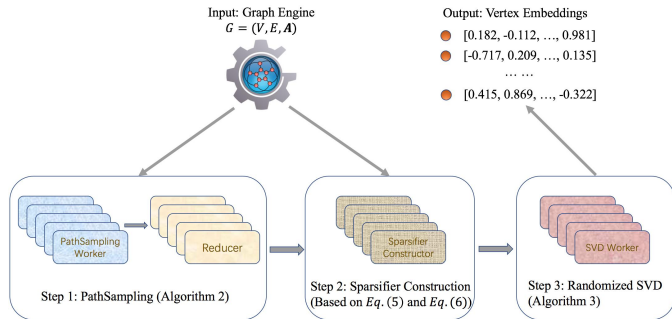


Figure 2: The System Design of NetSMF.

Revisit DeepWalk and NetMF

NetSMF: Network Embedding as Sparse Matrix Factorization

Experimental Results

Setup

Label Classification:

- ▶ BlogCatalog, PPI, Flickr, YouTube, OAG.
- ▶ Logistic Regression
- ▶ NetSMF ($T = 10$), NetMF ($T = 10$), DeepWalk, LINE.

Table 1: Statistics of Datasets.

Dataset	BlogCatalog	PPI	Flickr	YouTube	OAG
$ V $	10,312	3,890	80,513	1,138,499	67,768,244
$ E $	333,983	76,584	5,899,882	2,990,443	895,368,962
#Labels	39	50	195	47	19

Experimental Results

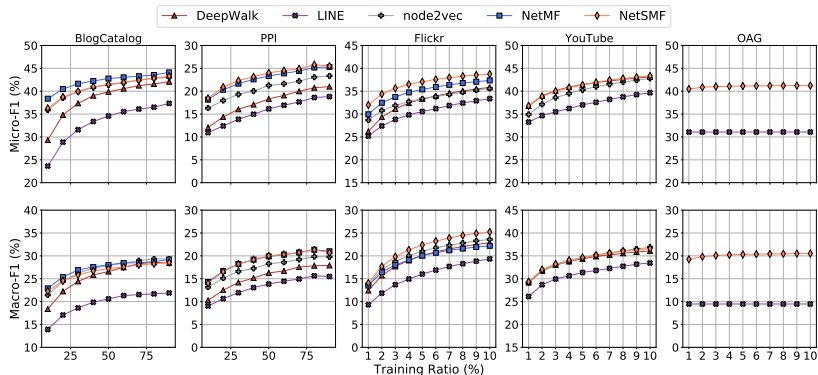


Figure 3: Predictive performance on varying the ratio of training data. The x-axis represents the ratio of labeled data (%), and the y-axis in the top and bottom rows denote the Micro-F1 and Macro-F1 scores respectively.

Running Time

Table 2: Running Time

	LINE	DeepWalk	node2vec	NetMF	NetSMF
BlogCatalog	40 mins	12 mins	56 mins	2 mins	13 mins
PPI	41 mins	4 mins	4 mins	16 secs	10 secs
Flickr	42 mins	2.2 hours	21 hours	2 hours	48 mins
YouTube	46 mins	1 day	4 days	×	4.1 hours
OAG	2.6 hours	–	–	×	24 hours

Conclusion and Future Work

We propose NetSMF, a scalable, efficient, and effective network embedding algorithm.

Future Work

- ▶ A distributed-memory implementation.
- ▶ Extension to directed, dynamic, heterogeneous graphs.

Thanks.

- ▶ Network Embedding as Matrix Factorization: Unifying DeepWalk, LINE, PTE, and node2vec (WSDM '18)
- ▶ NetSMF: Network Embedding as Sparse Matrix Factorization (WebConf '19)

Code for NetMF available at github.com/xptree/NetMF
Code for NetSMF available at github.com/xptree/NetSMF

Q&A

On the Large-dimensionality Assumption of [LG14]

Recall the objective of skip-gram model:

$$\min_{\mathbf{X}, \mathbf{Y}} \mathcal{L}(\mathbf{X}, \mathbf{Y})$$

where

$$\mathcal{L}(\mathbf{X}, \mathbf{Y}) = |\mathcal{D}| \sum_w \sum_c \left(\frac{\#(w, c)}{|\mathcal{D}|} \log g(\mathbf{x}_w^\top \mathbf{y}_c) + b \frac{\#(w)}{|\mathcal{D}|} \frac{\#(c)}{|\mathcal{D}|} \log g(-\mathbf{x}_w^\top \mathbf{y}_c) \right)$$

Theorem

For DeepWalk, when the length of random walk $L \rightarrow \infty$,

$$\frac{\#(w, c)}{|\mathcal{D}|} \xrightarrow{p} \frac{1}{2T} \sum_{r=1}^T \left(\frac{d_w}{\text{vol}(G)} (\mathbf{P}^r)_{w,c} + \frac{d_c}{\text{vol}(G)} (\mathbf{P}^r)_{c,w} \right).$$

$$\frac{\#(w)}{|\mathcal{D}|} \xrightarrow{p} \frac{d_w}{\text{vol}(G)} \text{ and } \frac{\#(c)}{|\mathcal{D}|} \xrightarrow{p} \frac{d_c}{\text{vol}(G)}.$$

NetSMF — Approximation Error

Denote $\mathbf{M} = \mathbf{D}^{-1}(\mathbf{D} - \mathbf{L})\mathbf{D}^{-1}$ in

$$\text{trunc_log}^\circ \left(\frac{\text{vol}(G)}{b} \mathbf{D}^{-1}(\mathbf{D} - \tilde{\mathbf{L}})\mathbf{D}^{-1} \right),$$

and $\tilde{\mathbf{M}}$ to be its sparsifier the we constructed.

Theorem

The singular value of $\tilde{\mathbf{M}} - \mathbf{M}$ satisfies

$$\sigma_i(\tilde{\mathbf{M}} - \mathbf{M}) \leq \frac{4\epsilon}{\sqrt{d_i d_{\min}}}, \forall i \in [n].$$

Theorem

Let $\|\cdot\|_F$ be the matrix Frobenius norm. Then

$$\left\| \text{trunc_log}^\circ \left(\frac{\text{vol}(G)}{b} \tilde{\mathbf{M}} \right) - \text{trunc_log}^\circ \left(\frac{\text{vol}(G)}{b} \mathbf{M} \right) \right\|_F \leq \frac{4\epsilon \text{vol}(G)}{b\sqrt{d_{\min}}} \sqrt{\sum_{i=1}^n \frac{1}{d_i}}.$$

Spectrally Similar

Definition

Suppose $G = (V, E, \mathbf{A})$ and $\tilde{G} = (V, \tilde{E}, \tilde{\mathbf{A}})$ are two weighted undirected networks. Let $\mathbf{L} = \mathbf{D}_G - \mathbf{A}$ and $\tilde{\mathbf{L}} = \mathbf{D}_{\tilde{G}} - \tilde{\mathbf{A}}$ be their Laplacian matrices, respectively. We define G and \tilde{G} are $(1 + \epsilon)$ -spectrally similar if

$$\forall \mathbf{x} \in \mathbb{R}^n, (1 - \epsilon) \cdot \mathbf{x}^\top \tilde{\mathbf{L}} \mathbf{x} \leq \mathbf{x}^\top \mathbf{L} \mathbf{x} \leq (1 + \epsilon) \cdot \mathbf{x}^\top \tilde{\mathbf{L}} \mathbf{x}.$$

NetSMF—Algorithm

Algorithm 1: NetSMF

Input : A social network $G = (V, E, \mathbf{A})$ which we want to learn network embedding;
The number of non-zeros M in the sparsifier; The dimension of embedding d .

Output: An embedding matrix of size $n \times d$, each row corresponding to a vertex.

```
1  $\tilde{G} \leftarrow (V, \emptyset, \tilde{\mathbf{A}} = \mathbf{0})$ 
  /* Create an empty network with  $E = \emptyset$  and  $\tilde{\mathbf{A}} = \mathbf{0}$ . */
2 for  $i \leftarrow 1$  to  $M$  do
3   Uniformly pick an edge  $e = (u, v) \in E$ 
4   Uniformly pick an integer  $r \in [T]$ 
5    $u', v', Z \leftarrow \text{PathSampling}(e, r)$ 
6   Add an edge  $(u', v', \frac{2rm}{MZ})$  to  $\tilde{G}$ 
   /* Parallel edges will be merged into one edge, with their weights
      summed up together. */
7 end
8 Compute  $\tilde{\mathbf{L}}$  to be the unnormalized graph Laplacian of  $\tilde{G}$ 
9 Compute  $\tilde{\mathbf{M}} = \mathbf{D}^{-1} (\mathbf{D} - \tilde{\mathbf{L}}) \mathbf{D}^{-1}$ 
10  $\mathbf{U}_d, \mathbf{\Sigma}_d, \mathbf{V}_d \leftarrow \text{RandomizedSVD}(\text{trunc\_log}^\circ \left( \frac{\text{vol}(G)}{b} \tilde{\mathbf{M}} \right), d)$ 
11 return  $\mathbf{U}_d \sqrt{\mathbf{\Sigma}_d}$  as network embeddings
```

Algorithm 2: PathSampling algorithm as described in [CCL⁺15].

- 1 **Procedure** PathSampling($e = (u, v), r$)
 - 2 Uniformly pick an integer $k \in [r]$
 - 3 Perform $(k - 1)$ -step random walk from u to u_0
 - 4 Perform $(r - k)$ -step random walk from v to u_r
 - 5 Keep track of $Z(\mathbf{p}) = \sum_{i=1}^r \frac{2}{A_{u_{i-1}, u_i}}$ along the length- r path \mathbf{p}
 between u_0 and u_r
 - 6 **return** $u_0, u_r, Z(\mathbf{p})$
-

Back

Randomized SVD

Algorithm 3: Randomized SVD on NetMF Matrix Sparsifier

```
1 Procedure RandomizedSVD( $A$ ,  $d$ )
2   Sampling Gaussian random matrix  $O$  //  $O \in \mathbb{R}^{n \times d}$ 
3   Compute sample matrix  $Y = A^T O = AO$  //  $Y \in \mathbb{R}^{n \times d}$ 
4   Orthonormalize  $Y$ 
5   Compute  $B = AY$  //  $B \in \mathbb{R}^{n \times d}$ 
6   Sample another Gaussian random matrix  $P$  //  $P \in \mathbb{R}^{d \times d}$ 
7   Compute sample matrix of  $Z = BP$  //  $Z \in \mathbb{R}^{n \times d}$ 
8   Orthonormalize  $Z$ 
9   Compute  $C = Z^T B$  //  $C \in \mathbb{R}^{d \times d}$ 
10  Run Jacobi SVD on  $C = U\Sigma V^T$ 
11  return  $ZU$ ,  $\Sigma$ ,  $YV$ 

/* Result matrices are of shape  $n \times d, d \times d, n \times d$  resp. */
```

Time and Space Complexity

Table 3: Time and Space Complexity of NetSMF.

	Time	Space
Step 1	$O(MT \log n)$ for weighted networks $O(MT)$ for unweighted networks	$O(M + n + m)$
Step 2	$O(M)$	$O(M + n)$
Step 3	$O(Md + nd^2 + d^3)$	$O(M + nd)$

References I

-  Dehua Cheng, Yu Cheng, Yan Liu, Richard Peng, and Shang-Hua Teng, *Spectral sparsification of random-walk matrix polynomials*, arXiv preprint arXiv:1502.03496 (2015).
-  Omer Levy and Yoav Goldberg, *Neural word embedding as implicit matrix factorization*, Advances in neural information processing systems, 2014, pp. 2177–2185.