

Network Embedding as Matrix Factorization: Unifying DeepWalk, LINE, PTE, and node2vec

Jiezhong Qiu

Tsinghua University

February 21, 2018

Joint work with Yuxiao Dong (MSR), Hao Ma (MSR),
Jian Li (IIIS, Tsinghua), Kuansan Wang (MSR),
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

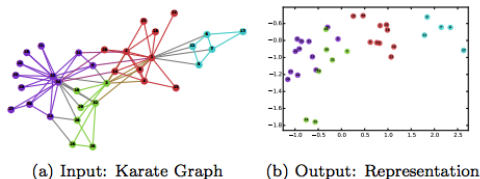
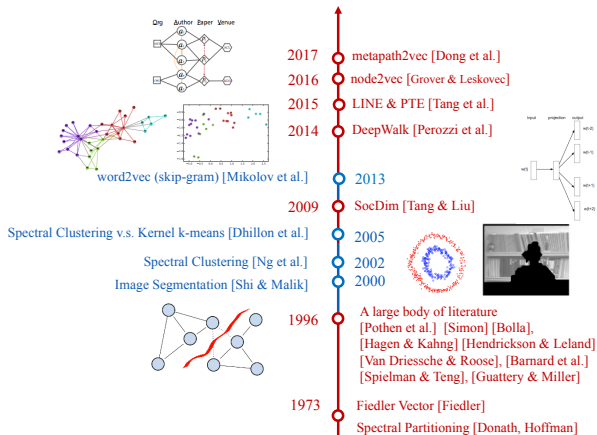


Figure 1: A toy example (Figure from DeepWalk).

History of Network Embedding



Contents

Preliminaries

Main Theoretic Results

Notations

DeepWalk (KDD'14)

LINE (WWW'15)

PTE (KDD'15)

node2vec (KDD'16)

NetMF

NetMF for a Small Window Size T

NetMF for a Large Window Size T

Experiments

Notations

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

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

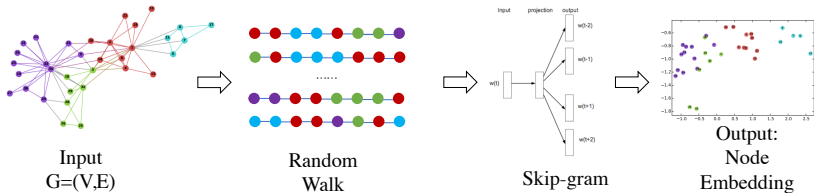
$$\mathbf{A}_{i,j} = \begin{cases} a_{i,j} > 0 & (i,j) \in E \\ 0 & (i,j) \notin E \end{cases}.$$

- ▶ 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 : $\text{vol}(G) = \sum_i \sum_j \mathbf{A}_{i,j}$.

Assumption

$G = (V, E)$ is connected, undirected, and not bipartite, which makes $P(w) = \frac{d_w}{\text{vol}(G)}$ a unique stationary distribution.

DeepWalk — Roadmap

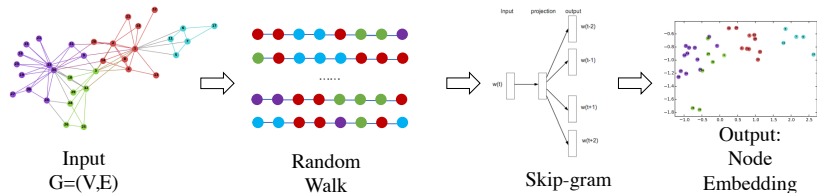


DeepWalk — a Two-step Algorithm

Algorithm 1: DeepWalk

```
1 for  $n = 1, 2, \dots, N$  do
2   Pick  $w_1^n$  according to a probability distribution  $P(w_1)$ ;
3   Generate a vertex sequence  $(w_1^n, \dots, w_L^n)$  of length  $L$  by a
   random walk on network  $G$ ;
4   for  $j = 1, 2, \dots, L - T$  do
5     for  $r = 1, \dots, T$  do
6       Add vertex-context pair  $(w_j^n, w_{j+r}^n)$  to multiset  $\mathcal{D}$ ;
7       Add vertex-context pair  $(w_{j+r}^n, w_j^n)$  to multiset  $\mathcal{D}$ ;
8 Run SGNS on  $\mathcal{D}$  with  $b$  negative samples.
```

DeepWalk — Roadmap



↓ 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

Skip-gram with Negative Sampling (SGNS)

- ▶ SGNS maintains a multiset \mathcal{D} which counts the occurrence of each word-context pair (w, c) .
- ▶ Objective:

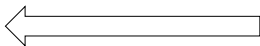
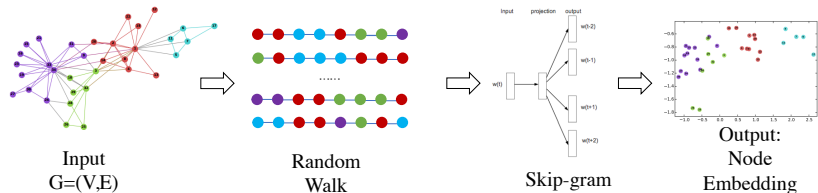
$$\mathcal{L} = \sum_w \sum_c \left(\#(w, c) \log g(\mathbf{x}_w^\top \mathbf{y}_c) + \frac{b \#(w) \#(c)}{|\mathcal{D}|} \log g(-\mathbf{x}_w^\top \mathbf{y}_c) \right),$$

where $\mathbf{x}_w, \mathbf{y}_c \in \mathbb{R}^d$, g is the sigmoid function, and b is the number of negative samples for SGNS.

- ▶ For sufficiently large dimensionality d , equivalent to factorizing PMI matrix (Levy & Goldberg, NIPS'14):

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

DeepWalk — Roadmap



Levy & Goldberg (NIPS 14)

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

b Number of negative samples

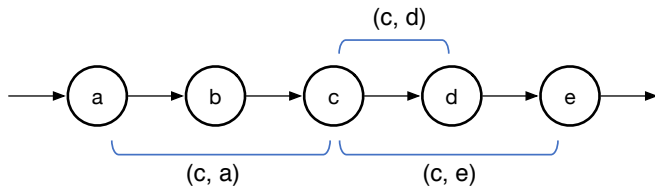
$\#(w, c)$ Co-occurrence of w and c

$\#(w)$ Occurrence of word w

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

$\#(c)$ Occurrence of context c

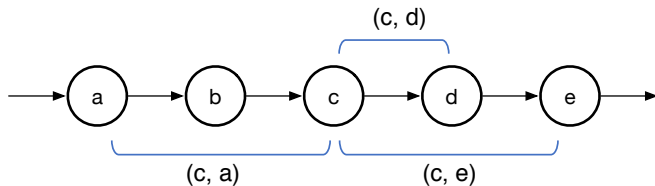
DeepWalk



Question

Suppose the multiset \mathcal{D} is constructed based on random walk on graph, can we interpret $\log \left(\frac{\#(w,c)|\mathcal{D}|}{b\#(w)\#(c)} \right)$ with graph theory terminologies?

DeepWalk



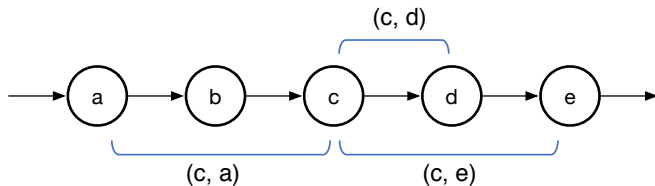
Question

Suppose the multiset \mathcal{D} is constructed based on random walk on graph, can we interpret $\log \left(\frac{\#(w,c)|\mathcal{D}|}{b\#(w)\#(c)} \right)$ with graph theory terminologies?

Challenge

We mix so many things together, i.e., direction and distance.

DeepWalk



Question

Suppose the multiset \mathcal{D} is constructed based on random walk on graph, can we interpret $\log \left(\frac{\#(w,c)|\mathcal{D}|}{b\#(w)\#(c)} \right)$ with graph theory terminologies?

Challenge

We mix so many things together, i.e., direction and distance.

Solution

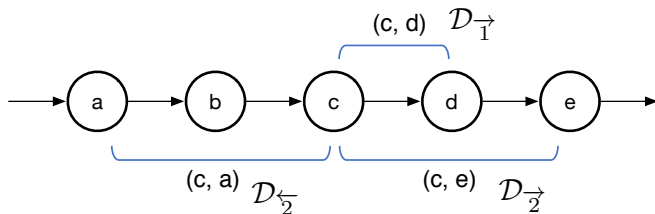
Let's distinguish them!

DeepWalk

Partition the multiset \mathcal{D} into several sub-multisets according to the way in which vertex and its context appear in a random walk sequence. More formally, for $r = 1, \dots, T$, we define

$$\mathcal{D}_{\vec{r}} = \{(w, c) : (w, c) \in \mathcal{D}, w = w_j^n, c = w_{j+r}^n\},$$

$$\mathcal{D}_{\overleftarrow{r}} = \{(w, c) : (w, c) \in \mathcal{D}, w = w_{j+r}^n, c = w_j^n\}.$$



DeepWalk as Implicit Matrix Factorization

Some observations

- ▶ Observation 1:

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

- ▶ Observation 2:

$$\frac{\#(w, c)}{|\mathcal{D}|} = \frac{1}{2T} \sum_{r=1}^T \left(\frac{\#(w, c)_{\vec{r}}}{|\mathcal{D}_{\vec{r}}|} + \frac{\#(w, c)_{\leftarrow r}}{|\mathcal{D}_{\leftarrow r}|} \right).$$

Sufficient to characterize $\frac{\#(w, c)_{\vec{r}}}{|\mathcal{D}_{\vec{r}}|}$ and $\frac{\#(w, c)_{\leftarrow r}}{|\mathcal{D}_{\leftarrow r}|}$.

DeepWalk — Theorems

Theorem

Denote $\mathbf{P} = \mathbf{D}^{-1}\mathbf{A}$, when the length of random walk $L \rightarrow \infty$,

$$\frac{\#(w, c)_{\vec{r}}}{|\mathcal{D}_{\vec{r}}|} \xrightarrow{p} \frac{d_w}{\text{vol}(G)} (\mathbf{P}^r)_{w,c} \quad \text{and} \quad \frac{\#(w, c)_{\overleftarrow{r}}}{|\mathcal{D}_{\overleftarrow{r}}|} \xrightarrow{p} \frac{d_c}{\text{vol}(G)} (\mathbf{P}^r)_{c,w}.$$

Theorem

When the length of random walk $L \rightarrow \infty$, we have

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

Theorem

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

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

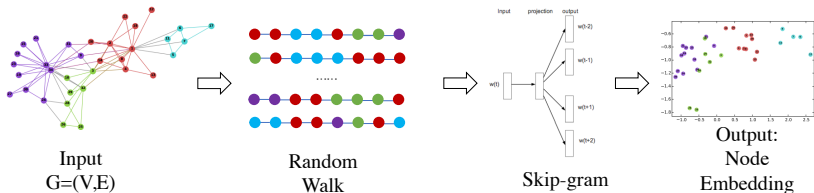
DeepWalk — Conclusion

Theorem

DeepWalk is asymptotically and implicitly factorizing

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

DeepWalk — Roadmap



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

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

\mathbf{A} Adjacency matrix

$$\text{vol}(G) = \sum_i \sum_j A_{i,j}$$

\mathbf{D} Degree matrix

b Number of negative samples

LINE

- ▶ Objective of LINE:

$$\mathcal{L} = \sum_{i=1}^{|V|} \sum_{j=1}^{|V|} \left(\mathbf{A}_{i,j} \log g(\mathbf{x}_i^\top \mathbf{y}_j) + \frac{bd_i d_j}{\text{vol}(G)} \log g(-\mathbf{x}_i^\top \mathbf{y}_j) \right).$$

- ▶ Align it with the Objective of SGNS:

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

- ▶ LINE is actually factorizing

$$\log \left(\frac{\text{vol}(G)}{b} \mathbf{D}^{-1} \mathbf{A} \mathbf{D}^{-1} \right)$$

- ▶ Recall DeepWalk's matrix form:

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

Observation LINE is a special case of DeepWalk ($T = 1$).

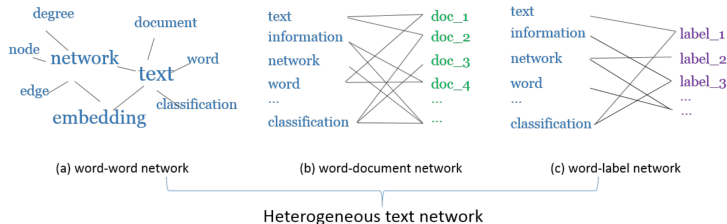


Figure 2: Heterogeneous Text Network.

- ▶ word-word network G_{ww} , $A_{ww} \in \mathbb{R}^{\#\text{word} \times \#\text{word}}$.
- ▶ document-word network G_{dw} , $A_{dw} \in \mathbb{R}^{\#\text{doc} \times \#\text{word}}$.
- ▶ label-word network G_{lw} , $A_{lw} \in \mathbb{R}^{\#\text{label} \times \#\text{word}}$.

PTE as Implicit Matrix Factorization

$$\log \left(\begin{bmatrix} \alpha \text{vol}(G_{ww})(D_{\text{row}}^{ww})^{-1} A_{ww} (D_{\text{col}}^{ww})^{-1} \\ \beta \text{vol}(G_{dw})(D_{\text{row}}^{dw})^{-1} A_{dw} (D_{\text{col}}^{dw})^{-1} \\ \gamma \text{vol}(G_{lw})(D_{\text{row}}^{lw})^{-1} A_{lw} (D_{\text{col}}^{lw})^{-1} \end{bmatrix} \right) - \log b,$$

- ▶ The matrix is of shape $(\#\text{word} + \#\text{doc} + \#\text{label}) \times \#\text{word}$.
- ▶ b is the number of negative samples in training.
- ▶ $\{\alpha, \beta, \gamma\}$ are hyper-parameters to balance the weights of the three networks. In PTE, $\{\alpha, \beta, \gamma\}$ satisfy

$$\alpha \text{vol}(G_{ww}) = \beta \text{vol}(G_{dw}) = \gamma \text{vol}(G_{lw})$$

node2vec — 2nd Order Random Walk

$$\underline{\mathbf{T}}_{u,v,w} = \begin{cases} \frac{1}{p} & (u,v) \in E, (v,w) \in E, u = w; \\ 1 & (u,v) \in E, (v,w) \in E, u \neq w, (w,u) \in E; \\ \frac{1}{q} & (u,v) \in E, (v,w) \in E, u \neq w, (w,u) \notin E; \\ 0 & \text{otherwise.} \end{cases}$$

$$\underline{\mathbf{P}}_{u,v,w} = \text{Prob}(w_{j+1} = u | w_j = v, w_{j-1} = w) = \frac{\underline{\mathbf{T}}_{u,v,w}}{\sum_u \underline{\mathbf{T}}_{u,v,w}}.$$

Stationary Distribution

$$\sum_w \underline{\mathbf{P}}_{u,v,w} \mathbf{X}_{v,w} = \mathbf{X}_{u,v}$$

Existence guaranteed by Perron-Frobenius theorem, but may not be unique.

node2vec as Implicit Matrix Factorization

Theorem

node2vec is asymptotically and implicitly factorizing a matrix whose entry at w -th row, c -th column is

$$\log \left(\frac{\frac{1}{2T} \sum_{r=1}^T (\sum_u \mathbf{X}_{w,u} \mathbf{P}_{c,w,u}^r + \sum_u \mathbf{X}_{c,u} \mathbf{P}_{w,c,u}^r)}{b (\sum_u \mathbf{X}_{w,u}) (\sum_u \mathbf{X}_{c,u})} \right)$$

Contents

Preliminaries

Main Theoretic Results

Notations

DeepWalk (KDD'14)

LINE (WWW'15)

PTE (KDD'15)

node2vec (KDD'16)

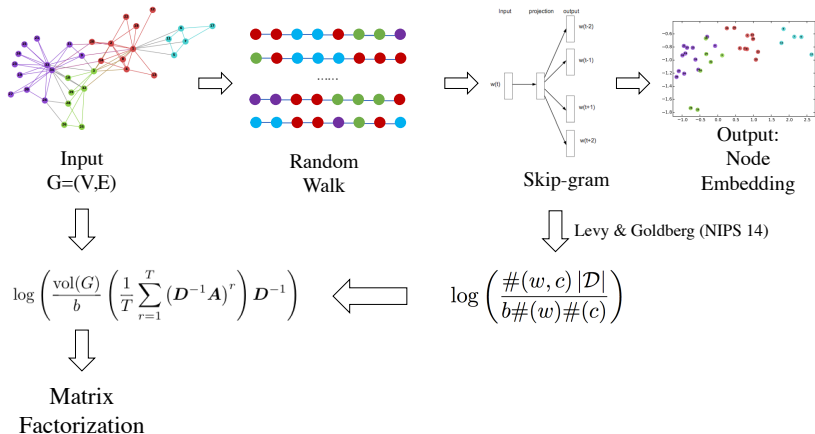
NetMF

NetMF for a Small Window Size T

NetMF for a Large Window Size T

Experiments

Roadmap



- ▶ Factorize the DeepWalk matrix:

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

- ▶ For numerical reason, we use truncated logarithm — $\tilde{\log}(x) = \log(\max(1, x))$

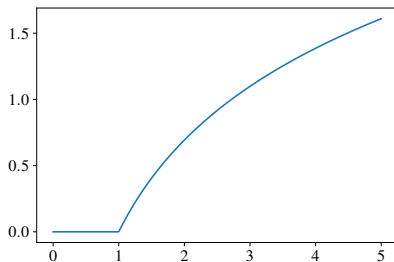


Figure 3: Truncated Logarithm

NetMF for a Small Window Size T

Algorithm 2: NetMF for a Small Window Size T

- 1 Compute $\mathbf{P}^1, \dots, \mathbf{P}^T$;
 - 2 Compute $\mathbf{M} = \frac{\text{vol}(G)}{bT} \left(\sum_{r=1}^T \mathbf{P}^r \right) \mathbf{D}^{-1}$;
 - 3 Compute $\mathbf{M}' = \max(\mathbf{M}, 1)$;
 - 4 Rank- d approximation by SVD: $\log \mathbf{M}' = \mathbf{U}_d \mathbf{\Sigma}_d \mathbf{V}_d^\top$;
 - 5 **return** $\mathbf{U}_d \sqrt{\mathbf{\Sigma}_d}$ as network embedding.
-

NetMF for a Large Window Size T — Observations

- ▶ We want to factorize

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

- ▶ We know the property of normalized graph Laplacian

$$\mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top$$

where $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_n)$ and $\forall \lambda_i \in [-1, 1]$.

$$\begin{aligned} \left(\frac{1}{T} \sum_{r=1}^T (\mathbf{D}^{-1} \mathbf{A})^r \right) \mathbf{D}^{-1} &= (\mathbf{D}^{-1/2}) \left(\frac{1}{T} \sum_{r=1}^T (\mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2})^r \right) (\mathbf{D}^{-1/2}) \\ &= (\mathbf{D}^{-1/2}) \left(\underbrace{\mathbf{U} \left(\frac{1}{T} \sum_{r=1}^T \mathbf{\Lambda}^r \right) \mathbf{U}^\top}_{\text{A polynomial}} \right) (\mathbf{D}^{-1/2}) \end{aligned}$$

NetMF for a Large Window Size T — Observations

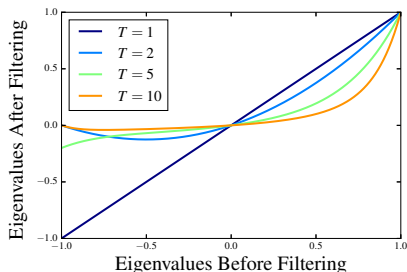


Figure 4: $f(\lambda) = \frac{1}{T} \sum_{r=1}^T \lambda^r$

Idea

This polynomial implicitly filters out negative eigenvalues and small positive eigenvalues, why not do it explicitly.

NetMF for a Large Window Size T — Algorithm

Algorithm 3: NetMF for a Large Window Size T

- 1 Eigen-decomposition $\mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2} \approx \mathbf{U}_h \mathbf{\Lambda}_h \mathbf{U}_h^\top$;
 - 2 Approximate \mathbf{M} with
$$\hat{\mathbf{M}} = \frac{\text{vol}(G)}{b} \mathbf{D}^{-1/2} \mathbf{U}_h \left(\frac{1}{T} \sum_{r=1}^T \mathbf{\Lambda}_h^r \right) \mathbf{U}_h^\top \mathbf{D}^{-1/2};$$
 - 3 Compute $\hat{\mathbf{M}}' = \max(\hat{\mathbf{M}}, 1)$;
 - 4 Rank- d approximation by SVD: $\log \hat{\mathbf{M}}' = \mathbf{U}_d \mathbf{\Sigma}_d \mathbf{V}_d^\top$;
 - 5 **return** $\mathbf{U}_d \sqrt{\mathbf{\Sigma}_d}$ as network embedding.
-

Setup

Label Classification:

- ▶ BlogCatalog, PPI, Wikipedia, Flickr
- ▶ Logistic Regression
- ▶ NetMF ($T = 1$) v.s. LINE
- ▶ NetMF ($T = 10$) v.s. DeepWalk

Table 1: Statistics of Datasets.

Dataset	BlogCatalog	PPI	Wikipedia	Flickr
$ V $	10,312	3,890	4,777	80,513
$ E $	333,983	76,584	184,812	5,899,882
#Labels	39	50	40	195

Experimental Results

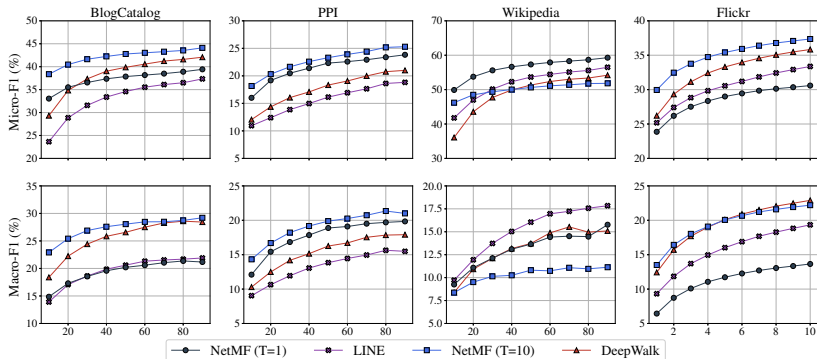


Figure 5: Predictive performance on varying the ratio of training data. The x-axis represents the ratio of labeled data (%), and the y-axis in the top and bottom rows denote the Micro-F1 and Macro-F1 scores respectively.

Conclusion

Table 2: The matrices that are implicitly approximated and factorized by DeepWalk, LINE, PTE, and node2vec.

Algorithm	Matrix
DeepWalk	$\log \left(\text{vol}(G) \left(\frac{1}{T} \sum_{r=1}^T (\mathbf{D}^{-1} \mathbf{A})^r \right) \mathbf{D}^{-1} \right) - \log b$
LINE	$\log \left(\text{vol}(G) \mathbf{D}^{-1} \mathbf{A} \mathbf{D}^{-1} \right) - \log b$
PTE	$\log \left(\begin{bmatrix} \alpha \text{vol}(G_{\text{ww}}) (\mathbf{D}_{\text{row}}^{\text{ww}})^{-1} \mathbf{A}_{\text{ww}} (\mathbf{D}_{\text{col}}^{\text{ww}})^{-1} \\ \beta \text{vol}(G_{\text{dw}}) (\mathbf{D}_{\text{row}}^{\text{dw}})^{-1} \mathbf{A}_{\text{dw}} (\mathbf{D}_{\text{col}}^{\text{dw}})^{-1} \\ \gamma \text{vol}(G_{\text{lw}}) (\mathbf{D}_{\text{row}}^{\text{lw}})^{-1} \mathbf{A}_{\text{lw}} (\mathbf{D}_{\text{col}}^{\text{lw}})^{-1} \end{bmatrix} \right) - \log b$
node2vec	$\log \left(\frac{\frac{1}{2T} \sum_{r=1}^T (\sum_u \mathbf{x}_{w,u} \mathbf{P}_{c,w,u}^r + \sum_u \mathbf{x}_{c,u} \mathbf{P}_{w,c,u}^r)}{(\sum_u \mathbf{x}_{w,u}) (\sum_u \mathbf{x}_{c,u})} \right) - \log b$

Thanks.

Standing on the shoulders of giants
— Isaac Newton

Code available at github.com/xptree/NetMF
Q&A

DeepWalk — Sketched Proof

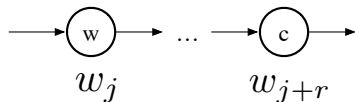
Theorem

Denote $\mathbf{P} = \mathbf{D}^{-1}\mathbf{A}$, when $L \rightarrow \infty$, we have

$$\frac{\#(w, c)_{\vec{r}}}{|\mathcal{D}_{\vec{r}}|} \xrightarrow{p} \frac{d_w}{\text{vol}(G)} (\mathbf{P}^r)_{w,c} \quad \text{and} \quad \frac{\#(w, c)_{\leftarrow r}}{|\mathcal{D}_{\leftarrow r}|} \xrightarrow{p} \frac{d_c}{\text{vol}(G)} (\mathbf{P}^r)_{c,w}.$$

Proof.

Consider the special case when $N = 1$, thus we only have one vertex sequence w_1, \dots, w_L generated by random walk. Let Y_j ($j = 1, \dots, L - T$) be the indicator function for event that $w_j = w$ and $w_{j+r} = c$



Proof (Con't)

Observation

- ▶ $\mathbb{E}[Y_j] = \text{Prob}(w_j = w, w_{j+r} = c) \rightarrow \frac{d_w}{\text{vol}(G)} (\mathbf{P}^r)_{w,c}$.
- ▶ $\frac{\#(w,c)_{\vec{r}}}{|\mathcal{D}_{\vec{r}}|} = \frac{1}{L-T} \sum_{j=1}^{L-T} Y_j$.
- ▶ $\text{Cov}(Y_i, Y_j) \rightarrow 0$ as $|i - j| \rightarrow \infty$.

Lemma

(S.N. Bernstein Law of Large Numbers) Let Y_1, Y_2, \dots be a sequence of random variables with finite expectation $\mathbb{E}[Y_j]$ and variance $\text{Var}(Y_j) < K$, $j \geq 1$, and covariances are s.t. $\text{Cov}(Y_i, Y_j) \rightarrow 0$ as $|i - j| \rightarrow \infty$. Then the law of large numbers (LLN) holds.

$$\frac{\#(w,c)_{\vec{r}}}{|\mathcal{D}_{\vec{r}}|} = \frac{1}{L-T} \sum_{j=1}^{L-T} Y_j \xrightarrow{p} \frac{1}{L-T} \sum_{j=1}^{L-T} \mathbb{E}(Y_j) \rightarrow \frac{d_w}{\text{vol}(G)} (\mathbf{P}^r)_{w,c}$$

Time Complexity

- ▶ Eigen-Decomposition (Implicitly Restarted Lanczos Method)
 $O(mhI + nh^2I + h^3I)$.
- ▶ Reconstruction $O(n^2h)$
- ▶ Element-wise logarithm $O(n^2)$.
- ▶ SVD (a naive implementation with eigen-decomposition):
 $O(n^2dI + nd^2I + d^3I)$.

Future Work

- ▶ Comprehend high-order cases, e.g., node2vec.

$$\log \left(\frac{\frac{1}{2T} \sum_{r=1}^T (\sum_u \mathbf{X}_{w,u} \mathbf{P}_{c,w,u}^r + \sum_u \mathbf{X}_{c,u} \mathbf{P}_{w,c,u}^r)}{b (\sum_u \mathbf{X}_{w,u}) (\sum_u \mathbf{X}_{c,u})} \right)$$

- ▶ Design scalable algorithm (e.g., using spectral sparsification of random-walk polynomials).

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

- ▶ Connection with graph convolutional networks (Kipf & Welling, ICLR'17).