



From word-embedding to node-embedding for learning over graphs

Fatemeh Sheikholeslami

DTC and Dept. of ECE, University of Minnesota

March 26, 2018

Taxonomy of learning over graphs



[1a] Perozzi, Bryan, Rami Al-Rfou, and Steven Skiena. "Deepwalk: Online learning of social representations." *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data mining*. ACM, New York, NY, Aug. 2014.

[1b] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S. Corrado, and Jeff Dean, "Distributed representations of words and phrases and their compositionality," In Advances in neural information processing systems, pp. 3111-3119, Lake Tahoe, USA, 2013.

[2] Thomas N. Kipf, Max Welling, Semi-Supervised Classification with Graph Convolutional Networks., ICLR, Toulon, France, April 2017.

[3] Aditya Grover and Jure Leskovec. "node2vec: Scalable feature learning for networks." *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge Discovery and Data mining*. ACM, San Francisco, CA, Aug. 2016.

Outline

- Taxonomy of learning over graphs
- Graph learning algorithms inspired by language models
 - ✓ Word2vec (Skip-gram)
 - ✓ DeepWalk
 - ✓ Node2vec
- Graph convolutional networks
- Simulation tests

Graph Representation

For machine learning tasks, the first step is to represent the network

- Such as adjacency matrix
- Feature extraction (node embedding)
 - Transform the adjacency matrix into a lower dimensional latent representation
 - Use deep learning techniques developed for language modeling



Language modeling

- Representation learning of words from documents
- The learned representations capture inherent structure

 \succ e.g., learn mapping $\phi: v \in V \to \mathbb{R}^{|V| \times d}$ such that

 $\|\phi(\text{rose}) - \phi(\text{daisy})\| < \|\phi(\text{rose}) - \phi(\text{tiger})\|$

- Recent approaches are based on word co-occurrence
- Word2vec
 - Continuous-bag-of-words
 - ✓ Skip-gram
 - ✓ GloVe

Word co-occurance maximization

- Consider source text "The quick brown fox jumps over the lazy dog."
- Skip-gram: with one word, predict surrounding words

$$> \text{ Max. the likelihood of context words}(w_{t-m}, \dots, w_{t-1}, w_{t+1}, \dots, w_{t+m}) \text{ given center word } w_t$$
$$\frac{1}{T} \sum_{t=1}^T \log P\Big((w_{t+m}, \dots, w_{t-1}, w_{t+1}, \dots, w_{t+m}) | w_t\Big) \simeq \frac{1}{T} \sum_{t=1}^T \sum_{j=-m,\dots,m; j \neq 0} \log P(w_{t+j} | w_t)$$



Window-size *m* is usually 5-10

Mikolov, Tomas, Kai Chen, Greg Corrado, and Jeffrey Dean. "Efficient estimation of word representations in vector space." *arXiv preprint arXiv:1301.3781* (2013).

Skipgram model

 \Box How to model $\log P(w_{t+j}|w_t)$?

Every word has 2 vectors for representation

- \succ \mathbf{v}_w : when w is the center word
- \succ \mathbf{u}_w : when w is the outside word (context word)

P

Softmax function

$$P(w_{t+j}|w_t) = \frac{\exp(\mathbf{u}_{w_{t+j}} \cdot \mathbf{v}_{w_t})}{\sum_{w'} \exp(\mathbf{u}_{w'} \cdot \mathbf{v}_{w_t})}$$



Neural network model

- Consider a vocabulary size of |V| = 10,000 and set d=300
- Shallow (one hidden layer) neural network model



- > Neural network input: one-hot vector corresponding to the "center word" of interest
- Maximization objective
- Gradient descent update rule

 $J(\boldsymbol{\theta}) = rac{1}{T} \sum_{t=1}^{T} J_t(\boldsymbol{\theta})$

 $\boldsymbol{\theta}^{new} = \boldsymbol{\theta}^{old} - \eta \nabla J(\boldsymbol{\theta}^{old})$

Final representation $\mathbf{\Phi} = \mathbf{U}^{ op} + \mathbf{V}$

Hidden Layer

✓ Stochastic gradient descent: randomly pick sample t $\theta^{new} = \theta^{old} - \eta \nabla J_t(\theta^{old})$

Word Vector

Algorithmic tweaks

- Word pairs and phrases
 - Find words that appear frequently together, e.g., "New York" and "Boston Globe".
 - "Phrase learning" is performed beforehand
- Subsampling of frequent words
 - Counter the imbalance between rare and frequent words
 - > Sample word w_i w/ frequency $f(w_i)$ w. p. $P(w_i)$

$$P(w_i) = 1 - \sqrt{\frac{t}{f(w_i)}}$$
 $t \simeq 10^{-5}$

x: 0.002537675 v: 1.0218046

Algorithmic tweaks- negative sampling

Typically extremely large vocabulary set ($|V| \simeq 10^5 - 10^7$)

Expensive to calculate the gradient of

$$P(w_o|w_c) = \frac{\exp(\mathbf{u}_{w_o} \cdot \mathbf{v}_{w_c})}{\sum_{w'} \exp(\mathbf{u}_{w'} \cdot \mathbf{v}_{w_c})}$$

> All |V|d weights in the second layer of the NN (**U**) are involved and updated!

- Recall that the output of the network is a one-hot vector.
- That is, one output is 1, and *all* of the other millions of output neurons are 0.
- Randomly select just a small number of "negative" words (say 5) to update the weights for.
 - "negative" words are the ones for which we want the network to output a 0
 - "positive" word is the word for which the network outputs 1
- Approximate the normalizer via few "negative samples"
- Probabilistic selection of negative samples

$$P(w_i) = rac{{f(w_i)}^{3/4}}{\sum_{j=0}^n \left({f(w_j)}^{3/4}
ight)}$$

Word analogies

- Nearest words to frog
 - 1. frogs
 - 2. toad
 - litoria
 - 4. leptodactylidae
 - 5. rana
 - 6. lizard
 - 7. eleutherodactylus

Test for linear relationships





litoria

leptodactylidae





rana

eleutherodactylus



Test results

Semantic-Syntactic Word Relationship test set

- Five types of semantic questions
- Nine types of syntactic questions

Type of relationship	Word	Pair 1	Word Pair 2		
Common capital city	Athens	Greece	Oslo	Norway	
All capital cities	Astana	Kazakhstan	Harare	Zimbabwe	
Currency	Angola	kwanza	Iran	rial	
City-in-state	Chicago	Illinois	Stockton	California	
Man-Woman	brother	sister	grandson	granddaughter	
Adjective to adverb	apparent	apparently	rapid	rapidly	
Opposite	possibly	impossibly	ethical	unethical	
Comparative	great	greater	tough	tougher	
Superlative	easy	easiest	lucky	luckiest	
Present Participle	think	thinking	read	reading	
Nationality adjective	Switzerland	Swiss	Cambodia	Cambodian	
Past tense	walking	walked	swimming	swam	
Plural nouns	mouse	mice	dollar	dollars	
Plural verbs	work	works	speak	speaks	

Accuracy on (subset of) the Semantic-Syntactic Word Relationship test set.

Dimensionality / Training words	24M	49M	98M	196M	391M	783M
50	13.4	15.7	18.6	19.1	22.5	23.2
100	19.4	23.1	27.8	28.7	33.4	32.2
300	23.2	29.2	35.3	38.6	43.7	45.9
600	24.0	30.1	36.5	40.8	46.6	50.4

Comparison of architectures using models trained on the same data, with d=640

Model	Semantic-Syntactic Wo	MSR Word Relatedness	
Architecture	Semantic Accuracy [%]	Syntactic Accuracy [%]	Test Set [20]
RNNLM	9	36	35
NNLM	23	53	47
CBOW	24	64	61
Skip-gram	55	59	56

Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S. Corrado, and Jeff Dean, "Distributed representations of words and phrases and their compositionality," In *Advances in neural information processing systems*, pp. 3111-3119, Lake Tahoe, USA, 2013. ¹²

From language modeling to graphs

Correspondence

- Words <--> Nodes
- Sentences <--> Node sequences
- Generating node sequences
 - Using random walks
 - short random walks = sentences

Connection

- Vertex frequency in random walks on scale-free graphs follows a power law.
- Words frequency in a natural language corpus follows a power law.

Deepwalk

1. Input: Graph

2.

3.

4.





-1.6

-0.5 0.0

-1.0

0.5

1.0

1.5 2.0

5. Representation-based inference such as classification, clustering, etc.

Bryan Perozzi, Rami Al-Rfou, and Steven Skiena, "Deepwalk: Online learning of social representations," *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data mining*. ACM, New York, NY, Aug. 2014. 14

Random walks

- Probability of going from node v to node x
 - > Common choice $\pi_{vx} = w_{vx}$





- 2nd order random walk [Grover et al. 2016)]
 - Consider a random walk that just traversed edge (t, v) and now resides at node v
 - > Probability of going to node x is set to $\pi_{vx} \sim \alpha_{p,q}(t,x) w_{vx}$

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

- p: Return parameter
- q: In-out parameter





Test results

Datasets

	BlogCatalog	PPI	Wikipedia
Ν	10,312	3,890	4,777
E	333,983	76,584	184,812
К	39	50	40

Performance vs. ratio of available labels



Roadmap

- Taxonomy of learning over graphs
- Graph learning algorithms inspired by NLP techniques
 - Skipgram
 - DeepWalk
 - Node2vec
- Graph convolutional networks
- Simulation tests

Learning over featured networks

◆ Consider a graph(known Adj.) of partly labeled N nodes w/ features $\mathbf{X} \in \mathbb{R}^{N \times F}$

Common practice

minimize $\mathcal{L} = \mathcal{L}_0 + \lambda \mathcal{L}_{reg}$, e.g. learn function f(.), w/ \mathcal{L}_0 = fitting term $\mathcal{L}_{reg} = \sum_{i,j} A_{ij} \|f(\mathbf{x}_i) - f(\mathbf{x}_j)\|^2$

Graph convolutional" approach

Use graph structure directly via a neural network model f(X, A)

Graph convolutional networks

Model

- > Neural network of depth *L*, with $\mathbf{X} \in \mathbb{R}^{N \times F}$ as input
- > $\mathbf{H}^{l} \in \mathbb{R}^{N \times D_{l}}$: Output matrix of the *l*-th layer

$$\mathbf{H}^{l+1} = \sigma(\tilde{\mathbf{D}}^{-1/2}\tilde{\mathbf{A}}\tilde{\mathbf{D}}^{-1/2}\mathbf{H}^{(l)}\mathbf{W}^{(l)})$$



- $\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{I}_N$, $\tilde{\mathbf{D}}_{ii} = \sum_i \tilde{A}_{ij}$, $\mathbf{H}^{(0)} = \mathbf{X}$, $\sigma(.)$: activation fun.
- Semi-supervised node classification
 - > Given one-hot labels Y of size $N_L \times K$ for N_L nodes

$$\min_{\mathbf{W}_0,\mathbf{W}_1} - \sum_{n=1}^{N_L} \sum_{k=1}^{K} Y_{nk} \ln Z_{nk} \qquad \mathbf{Z} = f(\mathbf{X}, \mathbf{A}) = \operatorname{softmax} \left(\check{\mathbf{A}} \operatorname{Relu}(\check{\mathbf{A}} \mathbf{X} \mathbf{W}^{(0)}) \mathbf{W}^{(1)} \right)$$

No explicit graph-based regularization

Thomas N. Kipf and Max Welling, "Semi-Supervised Classification with Graph Convolutional Networks," ICLR, Toulon, France, April 2017.

Test results

Datasets

Dataset	Туре	Nodes	Edges	Classes	Features	Label rate
Citeseer	Citation network	3,327	4,732	6	3,703	0.036
Cora	Citation network	2,708	5,429	7	1,433	0.052
Pubmed	Citation network	19,717	44,338	3	500	0.003
NELL	Knowledge graph	65,755	266,144	210	5,414	0.001

Classification accuracy (in percent)

Method	Citeseer	Cora	Pubmed	NELL
ManiReg [3]	60.1	59.5	70.7	21.8
SemiEmb [28]	59.6	59.0	71.1	26.7
LP [32]	45.3	68.0	63.0	26.5
DeepWalk [22]	43.2	67.2	65.3	58.1
ICA [18]	69.1	75.1	73.9	23.1
Planetoid* [29]	64.7 (26s)	75.7 (13s)	77.2 (25s)	61.9 (185s)
GCN (this paper)	70.3 (7s)	81.5 (4s)	79.0 (38s)	66.0 (48s)
GCN (rand. splits)	67.9 ± 0.5	80.1 ± 0.5	$\overline{78.9\pm0.7}$	58.4 ± 1.7

Semisupervised node embedding via GCN

Can be modified to serve as a feature extractor for nodes in a graph

 $\mathbf{Z} = f(\mathbf{X}, \mathbf{A}) = \operatorname{softmax} \left(\check{\mathbf{A}} \operatorname{Tanh} \left(\check{\mathbf{A}} \operatorname{Tanh} \left(\check{\mathbf{A}} \operatorname{Tanh} \left(\check{\mathbf{A}} \mathbf{W}^{(0)} \right) \mathbf{W}^{(1)} \right) \mathbf{W}^{(2)} \right)$

Train model via cross-entropy minimization on the labeled nodes

The hidden layer output is interpreted as the embedding vector



1.0 H

-1.0

-2

Effects of model depth

ResNet models for better performance on deep neural networks

$$\mathbf{H}^{l+1} = \sigma(\tilde{\mathbf{D}}^{-1/2}\tilde{\mathbf{A}}\tilde{\mathbf{D}}^{-1/2}\mathbf{H}^{(l)}\mathbf{W}^{(l)}) + \mathbf{H}^{(l)}$$

Performance



Summary

- Node embedding via language modeling tools
 - Skip-gram
 - DeepWalk
 - Node2vec
- Node embedding via GCN
 - Semi-supervised classification
 - Embedding via hidden layer outputs

