
Watch Your Step: Learning Node Embeddings via Graph Attention

Sami Abu-El-Haija*
Information Sciences Institute,
University of Southern California
haija@isi.edu

Bryan Perozzi
Google AI
New York City, NY
bperozzi@acm.org

Rami Al-Rfou
Google AI
Mountain View, CA
rmyeid@google.com

Alex Alemi
Google AI
Mountain View, CA
alemi@google.com

Abstract

Graph embedding methods represent nodes in a continuous vector space, preserving different types of relational information from the graph. There are many hyper-parameters to these methods (e.g. the length of a random walk) which have to be manually tuned for every graph. In this paper, we replace previously fixed hyper-parameters with trainable ones that we automatically learn via backpropagation. In particular, we propose a novel attention model on the power series of the transition matrix, which guides the random walk to optimize an upstream objective. Unlike previous approaches to attention models, the method that we propose utilizes attention parameters exclusively on the data traversal (e.g. on the random walk), and are not used by the model for inference. We experiment on link prediction tasks, as we aim to produce embeddings that best-preserve the graph structure, generalizing to unseen information. We improve state-of-the-art results on a comprehensive suite of real-world graph datasets including social, collaboration, and biological networks, where we observe that our graph attention model can reduce the error by 20% to 40%. We show that our automatically-learned attention parameters can vary significantly per graph, and correspond to the optimal choice of hyper-parameter if we manually tune existing methods.

1 Introduction

Unsupervised graph embedding methods seek to learn representations that encode the graph structure. These embeddings have demonstrated outstanding performance on a number of tasks including node classification [28, 13], knowledge-base completion [23], semi-supervised learning [35], and link prediction [2]. In general, as introduced by Perozzi et al [28], these methods operate in two discrete steps: First, they sample pair-wise relationships from the graph through random walks and counting node co-occurrences. Second, they train an embedding model e.g. using Skipgram or word2vec [24], to learn representations that encode pairwise node similarities.

While such methods have demonstrated positive results on a number of tasks, their performance can significantly vary based on the setting of their hyper-parameters. For example, [28] observed that the quality of learned representations is dependent on the length of the random walk (C). In practice, DeepWalk [28] and many of its extensions [e.g. 13] use word2vec implementations [24].

*Work has been done while Sami was at Google AI – formally, Google Research.

Accordingly, it has been revealed by [20] that the hyper-parameter C , referred to as *training window length* in word2vec [24], actually controls more than a fixed length of the random walk. Instead, it parameterizes a function, we term *the context distribution* and denote Q , which controls the probability of sampling a node-pair when visited within a specific distance². Implicitly, the choices of C and Q , create a weight mass on every node’s neighborhood. In general, the weight is higher on nearby nodes, but the specific form of the of the mass function is determined by the aforementioned hyper-parameters. In this work, we aim to replace these hyper-parameters with trainable parameters, so that they can be automatically learned for each graph. To do so, we pose graph embedding as end-to-end learning, where the (discrete) two steps of random walk co-occurrence sampling, followed by representation learning, are joint using a closed-form expectation over the graph adjacency matrix.

Our inspiration comes from the successful application of attention models in domains such as Natural Language Processing (NLP) [e.g. 4, 36], image recognition [25], and detecting rare events in videos [29]. To the best of our knowledge, the approach we propose is significantly different from the standard application of attention models. Instead of using attention parameters to guide the model where to look when making a prediction, we use attention parameters to guide our learning algorithm to focus on parts of the data that are most helpful for optimizing upstream objective.

We show mathematical equivalence between the context distribution and the co-efficients of power series of the transition matrix. This allows us to learn the context distribution by learning an attention model on the power series. The attention parameters “guide” the random walk, by allowing it to focus more on short- or long-term dependencies, as best suited for the graph, while optimizing an upstream objective. To the best of our knowledge, this work is the first application of attention methods to graph embedding.

Specifically, our contributions are the following:

1. We propose an extendible family of graph attention models that can learn arbitrary (e.g. non-monotonic) context distributions.
2. We show that the optimal choice of context distribution hyper-parameters for competing methods, found by manual tuning, agrees with our automatically-found attention parameters.
3. We evaluate on a number of challenging link prediction tasks comprised of real world datasets, including social, collaboration, and biological networks. Experiments show we substantially improve on our baselines, reducing link-prediction error by 20%-40%.

2 Preliminaries

2.1 Graph Embeddings

Given an unweighted graph $G = (V, E)$, its (sparse) adjacency matrix $\mathbf{A} \in \{0, 1\}^{|V| \times |V|}$ can be constructed according to $A_{vu} = \mathbb{1}[(v, u) \in E]$, where the indicator function $\mathbb{1}[\cdot]$ evaluates to 1 iff its boolean argument is true. In general, graph embedding methods minimize an objective:

$$\min_{\mathbf{Y}} \mathcal{L}(f(\mathbf{A}), g(\mathbf{Y})); \tag{1}$$

where $\mathbf{Y} \in \mathbb{R}^{|V| \times d}$ is a d -dimensional node embedding dictionary; $f : \mathbb{R}^{|V| \times |V|} \rightarrow \mathbb{R}^{|V| \times |V|}$ is a transformation of the adjacency matrix; $g : \mathbb{R}^{|V| \times d} \rightarrow \mathbb{R}^{|V| \times |V|}$ is a pairwise edge function; and $\mathcal{L} : \mathbb{R}^{|V| \times |V|} \times \mathbb{R}^{|V| \times |V|} \rightarrow \mathbb{R}$ is a loss function. For instance, a stochastic version of Singular Value Decomposition (SVD) is an embedding method, and can be casted into our framework by setting $f(\mathbf{A}) = \mathbf{A}$; decomposing \mathbf{Y} into two halves, the left- and right-embedding representations³ as $\mathbf{Y} = [\mathbf{L}|\mathbf{R}]$ with $\mathbf{L}, \mathbf{R} \in \mathbb{R}^{|V| \times \frac{d}{2}}$ then setting g to their outer product $g(\mathbf{Y}) = g([\mathbf{L}|\mathbf{R}]) = \mathbf{L} \times \mathbf{R}^\top$; and finally setting \mathcal{L} to the Frobenius norm of the error, yielding:

$$\min_{\mathbf{L}, \mathbf{R}} \|\mathbf{A} - \mathbf{L} \times \mathbf{R}^\top\|_F$$

²Specifically, rather than using C as constant and assuming all nodes visited within distance C are related, a desired context distance c_i is sampled from uniform ($c_i \sim \mathcal{U}\{1, C\}$) for each node pair i in training. If the node pair i was visited more than c_i -steps apart, it is not used for training. This was revealed to us by Levy et al [20], see their Section 3.1. Many DeepWalk-style methods inherited this, as they utilize word2vec implementation.

³Also known in NLP [24] as the “input” and “output” embedding representations.

2.2 Learning Embeddings via Random Walks

Introduced by [28], this family of methods [incl. 13, 18] induce random walks along E by starting from a random node $v_0 \in \text{sample}(V)$, and repeatedly sampling an edge to transition to next node as $v_{i+1} := \text{sample}(E[v_i])$, where $E[v_i]$ are the outgoing edges from v_i . The transition sequences $v_0 \rightarrow v_1 \rightarrow v_2 \rightarrow \dots$ (i.e. random walks) can then be passed to word2vec algorithm, which learns embeddings by stochastically taking every node along the sequence v_i , and the embedding representation of this **anchor** node v_i is brought closer to the embeddings of its next neighbors, $\{v_{i+1}, v_{i+2}, \dots, v_{i+c}\}$, the **context nodes**. In practice, the context window size c is sampled from a distribution e.g. uniform $\mathcal{U}\{1, C\}$ as explained in [20].

Let $\mathbf{D} \in \mathbb{R}^{|V| \times |V|}$ be the co-occurrence matrix from random walks, with each entry D_{vu} containing the number of times nodes v and u are co-visited within context distance $c \sim \mathcal{U}\{1, C\}$, in all simulated random walks. Embedding methods utilizing random walks, can also be viewed using the framework of Eq. (1). For example, to get Node2vec [13], we can set $f(\mathbf{A}) = \mathbf{D}$, set the edge function to the embeddings outer product $g(\mathbf{Y}) = \mathbf{Y} \times \mathbf{Y}^\top$, and set the loss function to negative log likelihood of softmax, yielding:

$$\min_{\mathbf{Y}} \left[\log Z - \sum_{v,u \in V} D_{vu} (Y_v^\top Y_u) \right], \quad (2)$$

where partition function $Z = \sum_{v,u} \exp(Y_v^\top Y_u)$ can be estimated with negative sampling [24, 13].

2.2.1 Graph Likelihood

A recently-proposed objective for learning embeddings is the *graph likelihood* [2]:

$$\prod_{v,u \in V} \sigma(g(\mathbf{Y})_{v,u})^{D_{vu}} (1 - \sigma(g(\mathbf{Y})_{v,u}))^{\mathbb{1}[(v,u) \notin E]}, \quad (3)$$

where $g(\mathbf{Y})_{v,u}$ is the output of the model evaluated at edge (v, u) , given node embeddings \mathbf{Y} ; the activation function $\sigma(\cdot)$ is the logistic; Maximizing the graph likelihood pushes the model score $g(\mathbf{Y})_{v,u}$ towards 1 if value D_{vu} is large and pushes it towards 0 if $(v, u) \notin E$.

In our work, we minimize the negative log of Equation 3, written in our matrix notation as:

$$\min_{\mathbf{Y}} \left\| -\mathbf{D} \circ \log(\sigma(g(\mathbf{Y}))) - \mathbb{1}[\mathbf{A} = 0] \circ \log(1 - \sigma(g(\mathbf{Y}))) \right\|_1, \quad (4)$$

which we minimize w.r.t node embeddings $\mathbf{Y} \in \mathbb{R}^{|V| \times d}$, where \circ is the Hadamard product; and the L1-norm $\|\cdot\|_1$ of a matrix is the sum of its entries. The entries of this matrix are positive because $0 < \sigma(\cdot) < 1$. Matrix $\mathbf{D} \in \mathbb{R}^{|V| \times |V|}$ can be created similar to [2], by counting node co-occurrences in simulated random walks.

2.3 Attention Models

We mention attention models that are most similar to ours [e.g. 25, 29, 33], where an attention function is employed to suggest positions within the input example that the classification function should *pay attention to, when making inference*. This function is used during the training phase in the forward pass and in the testing phase for prediction. The attention function and the classifier are jointly trained on an upstream objective e.g. cross entropy. In our case, the attention mechanism is only guides the learning procedure, and not used by the model for inference. Our mechanism suggests *parts of the data to focus on*, during training, as explained next.

3 Our Method

Following our general framework (Eq 1), we set $g(\mathbf{Y}) = g([\mathbf{L} \mid \mathbf{R}]) = \mathbf{L} \times \mathbf{R}^\top$ and $f(\mathbf{A}) = \mathbb{E}[\mathbf{D}]$, the expectation on co-occurrence matrix produced from simulated random walk. Using this closed form, we extend the NLGL loss (Eq. 4) to include attention parameters on the random walk sampling.

3.1 Expectation on the co-occurrence matrix: $\mathbb{E}[\mathbf{D}]$

Rather than obtaining \mathbf{D} by simulation of random walks and sampling co-occurrences, we formulate an expectation of this sampling, as $\mathbb{E}[\mathbf{D}]$. In general, this allows us to tune sampling parameters living inside of the random walk procedure including number of steps C .

Let \mathcal{T} be the transition matrix for a graph, which can be calculated by normalizing the rows of \mathbf{A} to sum to one. This can be written as:

$$\mathcal{T} = \text{diag}(\mathbf{A} \times \mathbf{1}_n)^{-1} \times \mathbf{A}. \quad (5)$$

Given an initial probability distribution $p^{(0)} \in \mathbb{R}^{|V|}$ of a random surfer, it is possible to find the distribution of the surfer after one step conditioned on $p^{(0)}$ as $p^{(1)} = p^{(0)T} \mathcal{T}$ and after k steps as $p^{(k)} = p^{(0)T} (\mathcal{T})^k$, where $(\mathcal{T})^k$ multiplies matrix \mathcal{T} with itself k -times. We are interested in an analytical expression for $\mathbb{E}[\mathbf{D}]$, the expectation over co-occurrence matrix produced by simulated random walks. A closed form expression for this matrix will allow us to perform end-to-end learning.

In practice, random walk methods based on DeepWalk [28] do not use C as a hard limit; instead, given walk sequence (v_1, v_2, \dots) , they sample $c_i \sim \mathcal{U}\{1, C\}$ separately for each anchor node v_i and potential context nodes, and only keep context nodes that are within c_i -steps of v_i . In expectation then, nodes $v_{i+1}, v_{i+2}, v_{i+3}, \dots$ will appear as context for anchor node v_i , respectively with probabilities $1, 1 - \frac{1}{C}, 1 - \frac{2}{C}, \dots$. We can write an expectation on $\mathbf{D} \in \mathbb{R}^{|V| \times |V|}$:

$$\mathbb{E} [\mathbf{D}^{\text{DEEPWALK}}; C] = \sum_{k=1}^C \Pr(c \geq k) \tilde{\mathbf{P}}^{(0)} (\mathcal{T})^k, \quad (6)$$

which is parametrized by the (discrete) walk length C ; where $\Pr(c \geq k)$ indicates the probability of node with distance k from anchor to be selected; and $\tilde{\mathbf{P}}^{(0)} \in \mathbb{R}^{|V| \times |V|}$ is a diagonal matrix (the initial positions matrix), with $\tilde{\mathbf{P}}_{vv}^{(0)}$ set to the number of walks starting at node v . Since $\Pr(c = k) = \frac{1}{C}$ for all $k = \{1, 2, \dots, C\}$, we can expand $\Pr(c \geq k) = \sum_{j=k}^C P(c = j)$, and re-write the expectation as:

$$\mathbb{E} [\mathbf{D}^{\text{DEEPWALK}}; C] = \tilde{\mathbf{P}}^{(0)} \sum_{k=1}^C \left[1 - \frac{k-1}{C} \right] (\mathcal{T})^k. \quad (7)$$

Eq. (3.1) is derived, step-by-step, in the Appendix. We are not concerned by the exact definition of the scalar coefficient, $\left[1 - \frac{k-1}{C} \right]$, but we note that the coefficient decreases with k .

Instead of keeping C a hyper-parameter, we want to analytically optimize it on an upstream objective. Further, we are interested to learn the co-efficients to $(\mathcal{T})^k$ instead of hand-engineering a formula.

As an aside, running the GloVe embedding algorithm [27] over the random walk sequences, in expectation, is equivalent to factorizing the co-occurrence matrix: $\mathbb{E} [\mathbf{D}^{\text{GloVe}}; C] = \tilde{\mathbf{P}}^{(0)} \sum_{k=1}^C \left[\frac{1}{k} \right] (\mathcal{T})^k$.

3.2 Learning the Context Distribution

We want to learn the co-efficients to $(\mathcal{T})^k$. Let the context distribution Q be a C -dimensional vector as $Q = (Q_1, Q_2, \dots, Q_C)$ with $Q_k \geq 0$ and $\sum_k Q_k = 1$. We assign co-efficient Q_k to $(\mathcal{T})^k$. Formally, our expectation on \mathbf{D} is parameterized with, and is differentiable w.r.t., Q :

$$\mathbb{E} [\mathbf{D}; Q_1, Q_2, \dots, Q_C] = \tilde{\mathbf{P}}^{(0)} \sum_{k=1}^C Q_k (\mathcal{T})^k = \tilde{\mathbf{P}}^{(0)} \mathbb{E}_{k \sim Q} [(\mathcal{T})^k], \quad (8)$$

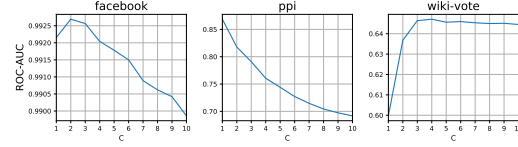
Training embeddings over random walk sequences, using word2vec or GloVe, respectively, are special cases of Equation 8, with Q fixed a priori as $Q_k = \left[1 - \frac{k-1}{C} \right]$ or $Q_k \propto \frac{1}{k}$.

3.3 Graph Attention Models

To learn Q automatically, we propose an attention model which guides the random surfer on “where to attend to” as a function of distance from the source node. Specifically, we define a *Graph Attention Model* as a process which models a node’s context distribution Q as the output of softmax:

$$(Q_1, Q_2, Q_3, \dots) = \text{softmax}((q_1, q_2, q_3, \dots)), \quad (9)$$

Dataset	$ V $	$ E $	nodes	edges
wiki-vote	7,066	103,663	users	votes
ego-Facebook	4,039	88,234	users	friendship
ca-AstroPh	17,903	197,031	researchers	co-authorship
ca-HepTh	8,638	24,827	researchers	co-authorship
PPI [31]	3,852	20,881	proteins	chemical interaction

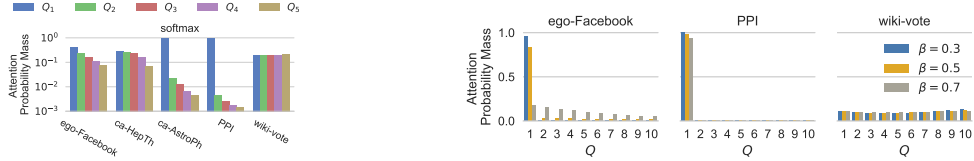


(a) Datasets used in our experiments.

(b) Test ROC-AUC as a function of C using node2vec.Figure 1: In 1a we present statistics of our datasets. In 1b, we motivate our work by showing the necessity of setting the parameter C for node2vec ($d=128$, each point is the average of 7 runs).

Dataset	dim	Adjacency Matrix			D by Simulation			Graph Attention (ours) (10)	Error Reduction
		Eigen Maps	SVD	DNGR	node2vec $C = 2$	node2vec $C = 5$	Asym Proj		
wiki-vote	64	61.3	86.0	59.8	64.4	63.6	91.7	93.8 ± 0.13	25.2%
	128	62.2	80.8	55.4	63.7	64.6	91.7	93.8 ± 0.05	25.2%
ego-Facebook	64	96.4	96.7	98.1	99.1	99.0	97.4	99.4 ± 0.10	33.3%
	128	95.4	94.5	98.4	99.3	99.2	97.3	99.5 ± 0.03	28.6%
ca-AstroPh	64	82.4	91.1	93.9	97.4	96.9	95.7	97.9 ± 0.21	19.2%
	128	82.9	92.4	96.8	97.7	97.5	95.7	98.1 ± 0.49	24.0%
ca-HepTh	64	80.2	79.3	86.8	90.6	91.8	90.3	93.6 ± 0.06	22.0%
	128	81.2	78.0	89.7	90.1	92.0	90.3	93.9 ± 0.05	23.8%
PPI	64	70.7	75.4	76.7	79.7	70.6	82.4	89.8 ± 1.05	43.5%
	128	73.7	71.2	76.9	81.8	74.4	83.9	91.0 ± 0.28	44.2%

Table 1: Results on Link Prediction Evaluation. Shown is the ROC-AUC.

(a) Learned Attention weights Q (log scale).(b) Q with varying the regularization β (linear scale).Figure 2: (a) shows learned attention weights Q , which agree with grid-search of node2vec (Figure 1b). (b) shows how varying β affects the learned Q . Note that distributions can quickly tail off to zero (ego-Facebook and PPI), while other graphs (wiki-vote) contain information across distant nodes.

where the variables q_k are trained via backpropagation, jointly while learning node embeddings. Our hypothesis is as follows. If we don't impose a specific formula on $Q = (Q_1, Q_2, \dots, Q_C)$, other than (regularized) softmax, then we can use very large values of C and allow every graph to learn its own form of Q with its preferred sparsity and own decay form. Should the graph structure require a small C , then the optimization would discover a left-skewed Q with all of probability mass on $\{Q_1, Q_2\}$ and $\sum_{k>2} Q_k \approx 0$. However, if according to the objective, a graph is more accurately encoded by making longer walks, then they can learn to use a large C (e.g. using uniform or even right-skewed Q distribution), focusing more attention on longer distance connections in the random walk.

To this end, we propose to train softmax attention model on the infinite power series of the transition matrix. We define an expectation on our proposed random walk matrix $\mathbf{D}^{\text{softmax}[\infty]}$ as⁴:

$$\mathbb{E} \left[\mathbf{D}^{\text{softmax}[\infty]}; q_1, q_2, q_3, \dots \right] = \tilde{\mathbf{P}}^{(0)} \lim_{C \rightarrow \infty} \sum_{k=1}^C \text{softmax}(q_1, q_2, q_3, \dots)_k (\mathcal{T})^k, \quad (10)$$

where q_1, q_2, \dots are jointly trained with the embeddings to minimize our objective.

3.4 Training Objective

The final training objective for the Softmax attention mechanism, coming from the NLGL Eq. (4),

$$\min_{\mathbf{L}, \mathbf{R}, \mathbf{q}} \beta \|\mathbf{q}\|_2^2 + \left\| -\mathbb{E}[\mathbf{D}; \mathbf{q}] \circ \log(\sigma(\mathbf{L} \times \mathbf{R}^\top)) - \mathbb{1}[\mathbf{A} = 0] \circ \log(1 - \sigma(\mathbf{L} \times \mathbf{R}^\top)) \right\|_1 \quad (11)$$

⁴We do not actually unroll the summation in Eq. (10) an infinite number of times. Our experiments show that unrolling it 10 or 20 times is sufficient to obtain state-of-the-art results.

is minimized w.r.t attention parameter vector $\mathbf{q} = (q_1, q_2, \dots)$ and node embeddings $\mathbf{L}, \mathbf{R} \in \mathbb{R}^{|V| \times \frac{d}{2}}$. Hyper-parameter $\beta \in \mathbb{R}$ applies L2 regularization on the attention parameters. We emphasize that our attention parameters \mathbf{q} live within the expectation over data \mathbf{D} , and are not part of the model (\mathbf{L}, \mathbf{R}) and are therefore not required for inference. The constraint $\sum_k Q_k = 1$, through the softmax activation, prevents $\mathbb{E}[\mathbf{D}^{\text{softmax}}]$ from collapsing into a trivial solution (zero matrix).

3.5 Algorithmic Complexity

The naive computation of $(\mathcal{T})^k$ requires k matrix multiplications and so is $\mathcal{O}(|V|^3 k)$. However, as most real-world adjacency matrices have an inherent low rank structure, a number of fast approximations to computing the random walk transition matrix raised to a power k have been proposed [e.g. 32]. Alternatively SVD can decompose \mathcal{T} as $\mathcal{T} = \mathcal{U}\Lambda\mathcal{V}^T$ and then the k^{th} power can be calculated by raising the diagonal matrix of singular values to k as $(\mathcal{T})^k = \mathcal{U}(\Lambda)^k\mathcal{V}^T$ since $\mathcal{V}^T\mathcal{U} = I$. Furthermore, the SVD can be approximated in time linear to the number of non-zero entries [14]. Therefore, we can calculate $(\mathcal{T})^k$ in $\mathcal{O}(|E|)$.

3.6 Extensions

As presented, our proposed method can learn the weights of the context distribution \mathcal{C} . However, we briefly note that such a model can be trivially extended to learn the weight of any other type of pair-wise node similarity (e.g. Personalized PageRank, Adamic-Adar, etc). In order to do this, we can extend the definition of the context Q with an additional dimension Q_{k+1} for the new type of similarity, and an additional element in the softmax q_{k+1} to learn a joint importance function.

4 Experiments

4.1 Link Prediction Experiments

We evaluate the quality of embeddings produced when random walks are augmented with attention, through experiments on link prediction [22]. Link prediction is a challenging task, with many real world applications in information retrieval, recommendation systems and social networks. As such, it has been used to study the properties of graph embeddings [28, 13]. Such an intrinsic evaluation emphasizes the structure-preserving properties of embedding.

Our experimental setup is designed to determine how well the embeddings produced by a method captures the topology of the graph. We measure this in the manner of [13]: remove a fraction (=50%) of graph edges, learn embeddings from the remaining edges, and measure how well the embeddings can recover those edges which have been removed. More formally, we split the graph edges E into two partitions of equal size E_{train} and E_{test} such that the training graph is connected. We also sample non existent edges $((u, v) \notin E)$ to make E_{train}^- and E_{test}^- . We use $(E_{\text{train}}, E_{\text{train}}^-)$ for training and model selection, and use $(E_{\text{test}}, E_{\text{test}}^-)$ to compute evaluation metrics.

Training: We train our models using TensorFlow, with PercentDelta optimizer [1]. For the results Table 1, we use $\beta = 0.5$, $C = 10$, and $\hat{\mathbf{P}}^{(0)} = \text{diag}(80)$, which corresponds to 80 walks per node. We analyze our model’s sensitivity in Section 4.2. To ensure repeatability of results, we have released our model and evaluation scripts.⁵

Datasets: Table 1a describes the datasets used in our experiments. Datasets available from SNAP <https://snap.stanford.edu/data>.

Baselines: We evaluate against many baselines. For all methods, we calculate $g(\mathbf{Y}) \in \mathbb{R}^{|V| \times |V|}$, and extract entries from $g(\mathbf{Y})$ corresponding to positive and negative test edges, then use them to compute ROC AUC. We compare against following baselines:

- **EigenMaps** [5]. Minimizes Euclidean distance of adjacent nodes of \mathbf{A} .
- **SVD**. Singular value decomposition of \mathbf{A} .
- **DNGR** [8]. Non-linear (i.e. deep) embedding of nodes, using an auto-encoder on \mathbf{A} . We use author’s code to learn the deep embeddings \mathbf{Y} and use for inference $g(\mathbf{Y}) = \mathbf{Y}\mathbf{Y}^T$.
- **node2vec** [13]. Simulates random walks and uses word2vec to learn node embeddings. Minimizes

⁵Available at <http://sami.haija.org/graph/attention.html>

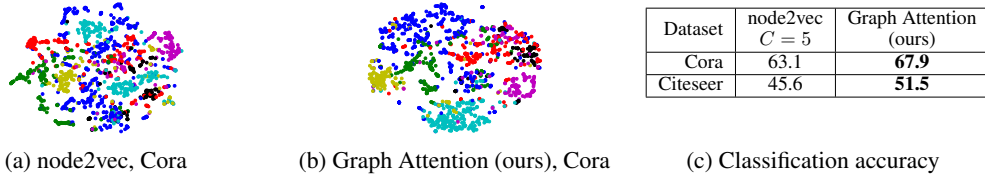


Figure 3: Node Classification. Fig. (a)/(b): t-SNE visualization of node embeddings for Cora dataset. We note that both methods are unsupervised, and we have colored the learned representations by node labels. Fig. (c) However, quantitatively, our embeddings achieves better separation.

objective in Eq. (2). For Table 1, we use author’s code to learn embeddings \mathbf{Y} then use $g(\mathbf{Y}) = \mathbf{Y}\mathbf{Y}^\top$. We run with $C = 2$ and $C = 5$.⁶

– **AsymProj** [2]. Learns edges as asymmetric projections in a deep embedding space, trained by maximizing the graph likelihood (Eq. 3). We take results from authors.

Results: Our results, summarized in Table 1, show that our proposed methods substantially outperform all baseline methods. Specifically, we see that the error is reduced by up to 45% over baseline methods which have fixed context definitions. This shows that by parameterizing the context distribution and allowing each graph to learn its own distribution, we can better preserve the graph structure (and thereby better predict missing edges).

Discussion: Figure 2a shows how the learned attention weights Q vary across datasets. Each dataset learns its own attention form, and the highest weights generally correspond to the highest weights when doing a grid search over C for node2vec (as in Figure 1b).

The hyper-parameter C determines the highest power of the transition matrix, and hence the maximum context size available to the attention model. We suggest using large values for C , since the attention weights can effectively use a subset of the transition matrix powers. For example, if a network needs only 2 hops to be accurately represented, then it is possible for the softmax attention model to learn $Q_3, Q_4, \dots \approx 0$. Figure 2b shows how varying the regularization term β allows the softmax attention model to “attend to” only what each dataset requires. We observe that for most graphs, the majority of the mass gets assigned to Q_1, Q_2 . This shows that shorter walks are more beneficial for most graphs. However, on wiki-vote, better embeddings are produced by paying attention to longer walks, as its softmax Q is uniform-like, with a slight right-skew.

4.2 Sensitivity Analysis

So far, we have removed two hyper-parameters, the maximum window size C , and the form of the context distribution \mathcal{U} . In exchange, we have introduced other hyper-parameters – specifically walk length (also C) and a regularization term β for the softmax attention model. Nonetheless, we show that our method is robust to various choices of these two. Figures 2a and 2b both show that the softmax attention weights drop to almost zero if the graph can be preserved using shorter walks, which is not possible with fixed-form distributions (e.g. \mathcal{U}).

Figure 4 examines this relationship in more detail for $d = 128$ dimensional embeddings, sweeping our hyper-parameters C and β , and comparing results to the best and worst node2vec embeddings for $C \in [1, 10]$. (Note that node2vec lines are horizontal, as they do not depend on β .) We observe that all the accuracy metrics are within 1% to 2%, when varying these hyper-parameters, and are all still well-above our baseline (which sample from a fixed-form context distribution).

4.3 Node Classification Experiments

We conduct node classification experiments, on two citation datasets, Cora and Citeseer, with the following statistics: Cora contains (2,708 nodes, 5,429 edges and $K = 7$ classes); and Citeseer contains (3,327 nodes, 4,732 edges and $K = 6$ classes). We learn embeddings from only the graph structure (nodes and edges), without observing node features nor labels during training. Figure 3 shows t-SNE visualization of the Cora dataset, comparing our method with node2vec [13]. For classification, we follow the data splits of [35]. We predict labels $\tilde{L} \in \mathbb{R}^{|V| \times K}$ as: $\tilde{L} = \exp(\alpha g(\mathbf{Y})) \times L_{\text{train}}$, where $L_{\text{train}} \in \{0, 1\}^{|V| \times K}$ contains rows of ones corresponding to

⁶We sweep C in Figure 1b, showing that there are no good default for C that works best across datasets.

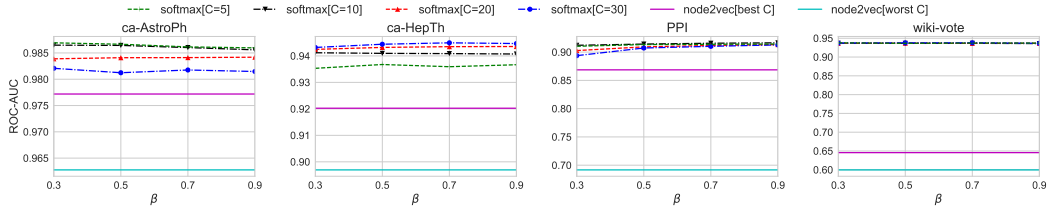


Figure 4: Sensitivity Analysis of softmax attention model. Our method is robust to choices of both β and C . We note that it consistently outperforms even an optimally set node2vec.

nodes in training set and zeros elsewhere. The scalar $\alpha \in \mathbb{R}$ is manually tuned on the validation set. The classification results, summarized in Table 3c, show that our model learns a better unsupervised representation than previous methods, that can then be used for supervised tasks. We do not compare against other semi-supervised methods that utilize node features during training and inference [incl. 35, 19], as our method is unsupervised.

Our classification prediction function contains one scalar parameter α . It can be thought of a “smooth” k-nearest-neighbors, as it takes a weighted average of known labels, where the weights are exponential of the dot-product similarity. Such a simple function should introduce no model bias.

5 Related Work

The field of learning on graphs has attracted much attention lately. Here we summarize two broad classes of algorithms, and point the reader to several recent reviews [6, 16] for more context.

The first class of algorithms are semi-supervised and concerned with predicting labels over a graph, its edges, and/or its nodes. Typically, these algorithms process a graph (nodes and edges) as well as per-node features. These include recent graph convolution methods [e.g. 26, 7, 3, 15, 33] with spectral variants [17, 11, 7], diffusion methods [e.g. 10, 12, 9], including ones trained until fixed-point convergence [30, 21] and semi-supervised node classification [35] with low-rank approximation of convolution [19]. We differ from these methods as (1) our algorithm is unsupervised (trained exclusively from the graph structure itself) without utilizing labels during training, and (2) we explicitly model the relationship between all node pairs.

The second class of algorithms consist of unsupervised graph embedding methods. Their primary goal is to preserve the graph structure, to create task independent representations. They explicitly model the relationship of all node pairs (e.g. as dot product of node embeddings). Some methods directly use the adjacency matrix [8, 34], and others incorporate higher order structure (e.g. from simulated random walks) [28, 13, 2]. Our work falls under this class of algorithms, where inference is a scoring function $V \times V \rightarrow \mathbb{R}$, trained to score positive edges higher than negative ones. Unlike existing methods, we do not specify a fixed context distribution apriori, whereas we push gradients through the random walk to those parameters, which we jointly train while learning the embeddings.

6 Conclusion

In this paper, we propose an attention mechanism for learning the context distribution used in graph embedding methods. We derive the closed-form expectation of the DeepWalk [28] co-occurrence statistics, showing an equivalence between the context distribution hyper-parameters, and the coefficients of the power series of the graph transition matrix. Then, we propose to replace the context hyper-parameters with trainable models, that we learn jointly with the embeddings on an objective that preserves the graph structure (the Negative Log Graph Likelihood, NLGL). Specifically, we propose Graph Attention Models, using a softmax to learn a free-form contexts distribution with a parameter for each type of context similarity (e.g. distance in a random walk).

We show significant improvements on link prediction and node classification over state-of-the-art baselines (that use a fixed-form context distribution), reducing error on link prediction and classification, respectively by up to 40% and 10%. In addition to improved performance (by learning distributions of arbitrary forms), our method can obviate the manual grid search over hyper-parameters:

walk length and form of context distribution, which can drastically fluctuate the quality of the learned embeddings and are different for every graph. On the datasets we consider, we show that our method is robust to its hyper-parameters, as described in Section 4.2. Our visualizations of converged attention weights convey to us that some graphs (e.g. voting graphs) can be better preserved by using longer walks, while other graphs (e.g. protein-protein interaction graphs) contain more information in short dependencies and require shorter walks.

We believe that our contribution in replacing these sampling hyperparameters with a learnable context distribution is general and can be applied to many domains and modeling techniques in graph representation learning.

References

- [1] S. Abu-El-Haija. Proportionate gradient updates with percentdelta. In *arXiv*, 2017.
- [2] S. Abu-El-Haija, B. Perozzi, and R. Al-Rfou. Learning edge representations via low-rank asymmetric projections. In *ACM International Conference on Information and Knowledge Management (CIKM)*, 2017.
- [3] J. Atwood and D. Towsley. Diffusion-convolutional neural networks. In *Advances in Neural Information Processing Systems (NIPS)*, 2016.
- [4] D. Bahdanau, K. Cho, and Y. Bengio. Neural machine translation by jointly learning to align and translate. In *International Conference on Learning Representations (ICLR)*, 2015.
- [5] M. Belkin and P. Niyogi. Laplacian eigenmaps for dimensionality reduction and data representation. In *Neural Computation*, 2003.
- [6] M. M. Bronstein, J. Bruna, Y. LeCun, A. Szlam, and P. Vandergheynst. Geometric deep learning: going beyond euclidean data. In *IEEE Signal Processing Magazine*, 2017.
- [7] J. Bruna, W. Zaremba, A. Szlam, and Y. LeCun. Spectral networks and deep locally connected networks on graphs. In *International Conference on Learning Representations*, 2013.
- [8] S. Cao, W. Lu, and Q. Xu. Deep neural networks for learning graph representations. In *Proceedings of the Association for the Advancement of Artificial Intelligence*, 2016.
- [9] L.-C. Chen, A. Schwing, A. Yuille, and R. Urtasun. Learning deep structured models. In *International Conference on Machine Learning*, 2015.
- [10] H. Dai, B. Dai, and L. Song. Discriminative embeddings of latent variable models for structured data. In *International Conference on Machine Learning*, 2016.
- [11] M. Defferrard, X. Bresson, and P. Vandergheynst. Convolutional neural networks on graphs with fast localized spectral filtering. In *Advances in Neural Information Processing Systems (NIPS)*, 2016.
- [12] D. Duvenaud, D. Maclaurin, J. Iparraguirre, R. Bombarell, T. Hirzel, A. Aspuru-Guzik, and R. Adams. Convolutional networks on graphs for learning molecular fingerprints. In *Advances in Neural Information Processing Systems (NIPS)*, 2015.
- [13] A. Grover and J. Leskovec. node2vec: Scalable feature learning for networks. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2016.
- [14] N. Halko, P. Martinsson, and J. A. Tropp. Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions. In *SIAM Review*, 2011.
- [15] W. Hamilton, R. Ying, and J. Leskovec. Inductive representation learning on large graphs. In *NIPS*, 2017.
- [16] W. L. Hamilton, R. Ying, and J. Leskovec. Representation learning on graphs: Methods and applications. In *IEEE Data Engineering Bulletin/IEEE Data Engineering Bulletin*, 2017.

- [17] M. Henaff, J. Bruna, and Y. LeCun. Deep convolutional networks on graph-structured data. In *arXiv:1506.05163*, 2015.
- [18] G. J. S. Ganguly, M. Gupta, V. Varma, and V. Pudi. Author2vec: Learning author representations by combining content and link information. In *Proceedings of the International Conference Companion on World Wide Web (WWW)*, WWW '16 Companion, 2016.
- [19] T. N. Kipf and M. Welling. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations (ICLR)*, 2017.
- [20] O. Levy, Y. Goldberg, and I. Dagan. Improving distributional similarity with lessons learned from word embeddings. In *TACL*, 2015.
- [21] Y. Li, D. Tarlow, M. Brockschmidt, and R. Zemel. Gated graph sequence neural networks. In *International Conference on Learning Representations*, 2016.
- [22] D. Liben-Nowell and J. Kleinberg. The link-prediction problem for social networks. In *Journal of American Society for Information Science and Technology*, 2007.
- [23] Y. Luo, Q. Wang, B. Wang, and L. Guo. Context-dependent knowledge graph embedding. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2015.
- [24] T. Mikolov, I. Sutskever, K. Chen, G. Corrado, and J. Dean. Distributed representations of words and phrases and their compositionality. In *Advances in Neural Information Processing Systems NIPS*. 2013.
- [25] V. Mnih, N. Heess, A. Graves, and k. kavukcuoglu. Recurrent models of visual attention. In *Advances in Neural Information Processing Systems (NIPS)*. 2014.
- [26] M. Niepert, M. Ahmed, and K. Kutzkov. Learning convolutional neural networks for graphs. In *International Conference on Machine Learning (ICML)*, 2016.
- [27] J. Pennington, R. Socher, and C. D. Manning. Glove: Global vectors for word representation. In *Conference on Empirical Methods in Natural Language Processing, EMNLP*, 2014.
- [28] B. Perozzi, R. Al-Rfou, and S. Skiena. Deepwalk: Online learning of social representations. In *Knowledge Discovery and Data Mining (KDD)*, 2014.
- [29] V. Ramanathan, J. Huang, S. Abu-El-Haija, A. Gorban, K. Murphy, and L. Fei-Fei. Detecting events and key actors in multi-person videos. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [30] F. Scarselli, M. Gori, A. Tsoi, M. Hagenbuchner, and G. Monfardini. The graph neural network model. In *IEEE Trans. on Neural Networks*, 2009.
- [31] C. Stark, B. Breitkreutz, T. Reguly, L. Boucher, A. Breitkreutz, and M. Tyers. Biogrid: A general repository for interaction datasets. In *Nucleic Acids Research*, 2006. URL <https://www.ncbi.nlm.nih.gov/pubmed/16381927>.
- [32] H. Tong, C. Faloutsos, and J.-Y. Pan. Fast random walk with restart and its applications. 2006.
- [33] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio. Graph attention networks. In *International Conference on Learning Representations (ICLR)*, 2018.
- [34] D. Wang, P. Cui, and W. Zhu. Structural deep network embedding. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2016.
- [35] Z. Yang, W. Cohen, and R. Salakhutdinov. Revisiting semi-supervised learning with graph embeddings. In *International Conference on Machine Learning (ICML)*, 2016.
- [36] Z. Yang, D. Yang1, C. Dyer, X. He, A. Smola, and E. Hovy. Hierarchical attention networks for document classification. In *Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)*, 2016.

7 Appendix

7.1 Step-by-step Derivation of Equation (3.1)

Let $x(k)$ be the position of random surfer at time k . Specifically, $x : \mathbb{Z}^+ \rightarrow V$. We assume a Markov chain: The value of $x(k)$ only depends on previous step: $x(k-1)$. To calculate the expectation $\mathbb{E}[\mathbf{D}]$, the square node-to-node co-occurrence matrix, we start by calculating one entry at a time: $\mathbb{E}[D_{uv}]$, the expected number of times that u is selected in v 's context. Let $W_v(k)$ be the *context set* that gets sampled if v is visited at the k^{th} step. Concretely, if $x(k) = v$, and the random walker continues the sequence, $x(k+1) = v'_1$ then $x(k+2) = v'_2$ then $x(k+3) = v'_3 \dots$, the context set of DeepWalk can be defined as $W_v(k) = \{v'_1, v'_2 \dots, v'_c\}$, where $c \sim \mathcal{U}\{1, C\}$. We would like to count the event $u \in W_v(k)$ for every $k \in \{1, 2, \dots, C\}$.

Using Markov Chain, we can write:

$$\begin{aligned} \Pr(x(i+k) = u \mid x(i) = v) &= \Pr(x(k) = u \mid x(0) = v) \\ &= (\mathcal{T}^k)_{uv} \end{aligned} \quad (12)$$

Now, if node u was visited k steps after node v , then the probability of it being sampled is given by:

$$\Pr(u \in W_v(k) \mid x(k) = u, x(0) = v). \quad (13)$$

In case of DeepWalk [28], probability above equals:

$$\Pr(c \geq k \mid x(k) = u, x(0) = v) \text{ where } c \sim \mathcal{U}\{1, C\}, \quad (14)$$

and event $k \leq c$ is independent of the condition $(x(k) = u \cap x(0) = v)$. Further, event $k \leq c$ can be partitioned and Eq. (14) can be written as

$$\Pr(c = k \cup c = k+1 \cup \dots \cup c = C) \quad (15)$$

$$= \sum_{j=k}^C \Pr(c = j) \quad (16)$$

$$= (C - k + 1) \left(\frac{1}{C}\right) = 1 - \frac{k-1}{C}, \quad (17)$$

where second line is trivial since the events $c = j$ are disjoint. We can now use Bayes' rule to derive the probability of u being visited k steps after v and being selected in v 's sampled context, as:

$$\begin{aligned} \Pr(u \in W_v(k), x(k) = u \mid x(0) = v) \\ &= \Pr(u \in W_v(k) \mid x(k) = u, x(0) = v) \Pr(x(k) = u \mid x(0) = v) \\ &= \left(1 - \frac{k-1}{C}\right) (\mathcal{T}^k)_{uv} \end{aligned} \quad (18)$$

Now, let E_{vku} be the event that a walker visits v and after k steps, visits u and selects it part of its context. This event happens with the probability indicated in Equation 18. Concretely,

$$\mathbb{E}[E_{vku} \mid x(0) = v] = \left(1 - \frac{k-1}{C}\right) (\mathcal{T}^k)_{uv}. \quad (19)$$

Let E_{v*u} count the events $\{E_{vku} : k \in [1, C]\}$, then:

$$\mathbb{E}[E_{v*u} \mid x(0) = v] = \mathbb{E}\left[\sum_{k=1}^C E_{vku} \mid x(0) = v\right] \quad (20)$$

$$= \sum_{k=1}^C \mathbb{E}[E_{vku} \mid x(0) = v] = \sum_{k=1}^C \left(1 - \frac{k-1}{C}\right) (\mathcal{T}^k)_{uv}. \quad (21)$$

Suppose we run DeepWalk, starting m random walks from each node v , then the expected number of times that u is present in the context of v is given by:

$$\mathbb{E}[D_{uv}^{\text{DEEPWALK}}] = m \mathbb{E}[E_{v*u} \mid x(0) = v] = m \sum_{k=1}^C \left(1 - \frac{k-1}{C}\right) (\mathcal{T}^k)_{uv}.$$

Finally, we can write down the expectation over the square matrix \mathbf{D} :

$$\begin{aligned} \mathbb{E} [\mathbf{D}^{\text{DEEPWALK}}] &= \text{diag}(m, m, \dots, m) \sum_{k=1}^C \left(1 - \frac{k-1}{C}\right) (\mathcal{T}^k) \\ &= \text{Equation (3.1)} \end{aligned}$$

□

7.2 Choice of $\tilde{\mathbf{P}}^{(0)}$

The github code of DeepWalk and node2vec start a fixed number (m) walks from every graph node $v \in V$. For node2vec, m defaults to 10, see num-walks flag in <https://github.com/aditya-grover/node2vec/blob/master/src/main.py>. Therefore, in our experiments, we set $\tilde{\mathbf{P}}^{(0)} := \text{diag}(m, m, \dots, m)$. This initial condition yields D_{vu} to be the expected number of times that u is visited if we started m walks from v . There can be other reasonable choices. Nonetheless, we use what worked well in practice for [13, 28]. We leave the search for a better $\tilde{\mathbf{P}}^{(0)}$ as future work.

7.3 Depiction of Learned Context Distribution

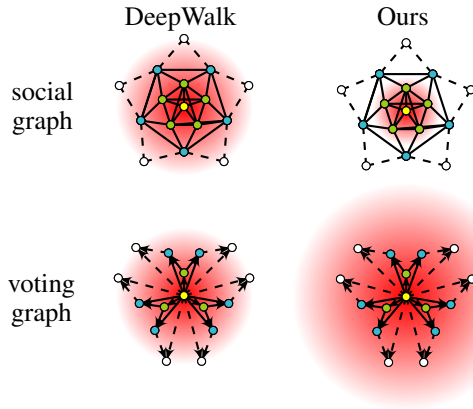


Figure: Depiction of how our model assigns context distributions (shaded red) compared to earlier work. We depict the graph from the perspective of anchor node (yellow). Given a social graph (top), where friends of friends are usually friends, our algorithm learns a leftskewed distribution. Given a voting graph (bottom), with general transitivity: $a \rightarrow b \rightarrow c \implies a \rightarrow c$, it learns a long-tail distribution. Earlier methods (e.g. DeepWalk) use word2vec, which internally uses a linear decay context distribution, treating all graphs the same.