
N-GCN: Multi-scale Graph Convolution for Semi-supervised Node Classification

Sami Abu-El-Haija*
Univ. of Southern California
Los Angeles, CA, USA
sami@haija.org

Amol Kapoor
Google Research
New York, NY, USA
ajkapoor@google.com

Bryan Perozzi
Google Research
New York, NY, USA
bperozzi@acm.org

Joonseok Lee
Google Research
Mountain View, CA, USA
joonseok@google.com

Abstract

Graph Convolutional Networks (GCNs) have shown significant improvements in semi-supervised learning on graph-structured data. Concurrently, unsupervised learning of graph embeddings has benefited from the information contained in random walks. In this paper, we propose a model: Network of GCNs (N-GCN), which marries these two lines of work. At its core, N-GCN trains multiple instances of GCNs over node pairs discovered at different distances in random walks, and learns a combination of the instance outputs which optimizes the classification objective. Our experiments show that our proposed N-GCN model improves state-of-the-art baselines on all of the challenging node classification tasks we consider: Cora, Citeseer, Pubmed, and PPI. In addition, our proposed method has other desirable properties, including generalization to recently proposed semi-supervised learning methods such as GraphSAGE, allowing us to propose N-SAGE, and resilience to adversarial input perturbations.¹

1 INTRODUCTION

Semi-supervised learning on graphs is important in many real-world applications, where the goal is to recover labels for all nodes given only a fraction of labeled ones. Some applications include social networks, where one wishes to predict user interests, or in health care, where one wishes to predict whether a patient should be screened for cancer. In many such cases, collecting node labels can be prohibitive. However, edges between nodes

can be easier to obtain, either using an explicit graph (e.g., social network) or implicitly by calculating pairwise similarities (e.g., using a patient-patient similarity kernel, Merdan et al., 2017).

Convolutional Neural Networks (LeCun et al., 1998) learn location-invariant hierarchical filters, enabling significant improvements on Computer Vision tasks (Krizhevsky et al., 2012; Szegedy et al., 2015; He et al., 2016). This success has motivated researchers (Bruna et al., 2014) to extend convolutions from spatial (i.e., regular lattice) domains to graph-structured (i.e., irregular) domains, yielding a class of algorithms known as Graph Convolutional Networks (GCNs).

Formally, we are interested in semi-supervised learning where we are given a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ with $N = |\mathcal{V}|$ nodes; adjacency matrix \mathbf{A} ; and matrix $\mathbf{X} \in \mathbb{R}^{N \times F}$ of node features, where F is the feature dimensionality. Labels for only a subset of nodes $\mathcal{V}_L \subset \mathcal{V}$ are observed. In general, $|\mathcal{V}_L| \ll |\mathcal{V}|$. Our goal is to recover labels for all unlabeled nodes $\mathcal{V}_U = \mathcal{V} - \mathcal{V}_L$, using the feature matrix \mathbf{X} , the known labels for nodes in \mathcal{V}_L , and the graph \mathcal{G} . In this setting, one treats the graph as the “unsupervised” and labels of \mathcal{V}_L as the “supervised” portions of the data.

Depicted in Fig. 1, our model for semi-supervised node classification builds on the GCN module proposed by Kipf and Welling (2017), which operates on the normalized adjacency matrix $\hat{\mathbf{A}}$, as in $\text{GCN}(\hat{\mathbf{A}})$, where $\hat{\mathbf{A}} = \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}}$, and \mathbf{D} is diagonal matrix of node degrees. Our proposed extension of GCNs is inspired by the recent advancements in random walk based graph embeddings (e.g. Perozzi et al., 2014; Grover and Leskovec, 2016; Abu-El-Haija et al., 2018). We make a Network of GCN modules (N-GCN), feeding each module a different power of $\hat{\mathbf{A}}$, as in $\{\text{GCN}(\hat{\mathbf{A}}^0), \text{GCN}(\hat{\mathbf{A}}^1), \text{GCN}(\hat{\mathbf{A}}^2), \dots\}$. The k -th power contains statistics from the k -th step of a random walk on the graph. Therefore, our N-GCN model is able to combine information from various step-sizes

This work was done at Google Research.

¹Source code: <https://github.com/samihaija/mixhop>

(i.e. graph scales). We then combine the output of all GCN modules into a classification sub-network, and we jointly train all GCN modules and the classification sub-network on the upstream objective for semi-supervised node classification. Weights of the classification sub-network give us insight on how the N-GCN model works. For instance, in the presence of input perturbations, we observe that the classification sub-network weights shift towards GCN modules utilizing higher powers of the adjacency matrix, effectively widening the “receptive field” of the (spectral) convolutional filters. We achieve state-of-the-art on several semi-supervised graph learning tasks, showing that explicit random walks enhance the representational power of vanilla GCN’s.

The rest of this paper is organized as follows. Section 2 reviews background work that provides the foundation for this paper. In Section 3, we describe our proposed method, followed by experimental evaluation in Section 4. We compare our work with recent closely-related methods in Section 5. Finally, we conclude with our contributions and future work in Section 6.

2 BACKGROUND

2.1 Semi-Supervised Node Classification

Traditional label propagation algorithms (Weston et al., 2012; Belkin et al., 2006) learn a model that transforms node features into node labels and uses the graph to add a regularizer term:

$$\mathcal{L}_{\text{label,prop}} = \mathcal{L}_{\text{cls}}(f(\mathbf{X}), \mathcal{V}_L) + \lambda f(\mathbf{X})^\top \Delta f(\mathbf{X}). \quad (1)$$

The first term \mathcal{L}_{cls} (classification loss) trains the model $f : \mathbb{R}^{N \times F} \rightarrow \mathbb{R}^{N \times C}$ to predict the known labels \mathcal{V}_L , where F is the feature dimensionality and C is the number of classes. The second term is the graph-based regularizer, ensuring that connected nodes have a similar model output, with Δ being the graph Laplacian and $\lambda \in \mathbb{R}$ is the regularization coefficient hyperparameter.

2.2 Graph Convolutional Networks

Graph Convolution (Bruna et al., 2014) generalizes convolution from Euclidean domains to graph-structured data. Convolving a “filter” over a signal on graph nodes can be calculated by transforming both the filter and the signal to the Fourier domain, multiplying them, and then transforming the result back into the discrete domain. The signal transform is achieved by multiplying with the eigenvectors of the graph Laplacian. The transformation requires a quadratic eigendecomposition of the symmetric Laplacian; however, the low-rank approximation

of the eigendecomposition can be calculated using truncated Chebyshev polynomials (Hammond et al., 2011). For instance, Kipf and Welling (2017) calculates a rank-1 approximation of the decomposition. They propose a multi-layer Graph Convolutional Networks (GCNs) for semi-supervised graph learning. Every layer computes the transformation:

$$\mathbf{H}^{(l+1)} = \sigma \left(\hat{\mathbf{A}} \mathbf{H}^{(l)} \mathbf{W}^{(l)} \right), \quad (2)$$

where $\mathbf{H}^{(l)} \in \mathbb{R}^{N \times d_l}$ is the input activation matrix to the l -th hidden layer with row $\mathbf{H}_i^{(l)}$ containing a d_l -dimensional feature vector for vertex $i \in \mathcal{V}$, and $\mathbf{W}^{(l)} \in \mathbb{R}^{d_l \times d_{l+1}}$ is the layer’s trainable weights. The first hidden layer $\mathbf{H}^{(0)}$ is set to the input features \mathbf{X} . A softmax on the last layer is used to classify labels. All layers use the same “normalized adjacency” $\hat{\mathbf{A}}$, obtained by the “renormalization trick” utilized by Kipf and Welling (2017), as $\hat{\mathbf{A}} = \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}}$,²

Eq. (2) is a first order approximation of convolving filter $\mathbf{W}^{(l)}$ over signal $\mathbf{H}^{(l)}$ (Hammond et al., 2011; Kipf and Welling, 2017). The left-multiplication with $\hat{\mathbf{A}}$ averages node features with their direct neighbors; this signal is then passed through a non-linearity function $\sigma(\cdot)$ (e.g., $\text{ReLU}(z) = \max(0, z)$). Successive layers effectively *diffuse* signals from nodes to neighbors.

Two-layer GCN model can be defined in terms of vertex features \mathbf{X} and normalized adjacency $\hat{\mathbf{A}}$ as:

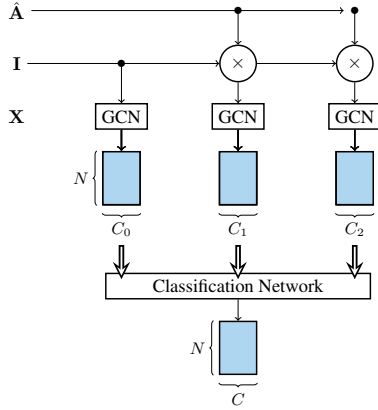
$$\text{GCN}_{2\text{-layer}}(\hat{\mathbf{A}}, \mathbf{X}; \theta) = \text{softmax} \left(\hat{\mathbf{A}} \sigma(\hat{\mathbf{A}} \mathbf{X} \mathbf{W}^{(0)}) \mathbf{W}^{(1)} \right),$$

where the GCN parameters $\theta = \{\mathbf{W}^{(0)}, \mathbf{W}^{(1)}\}$ are trained to minimize the cross-entropy error over labeled examples. The output of the GCN model is a matrix $\mathbb{R}^{N \times C}$, where N is the number of nodes and C is the number of labels. Each row contains the label scores for one node, assuming there are C classes.

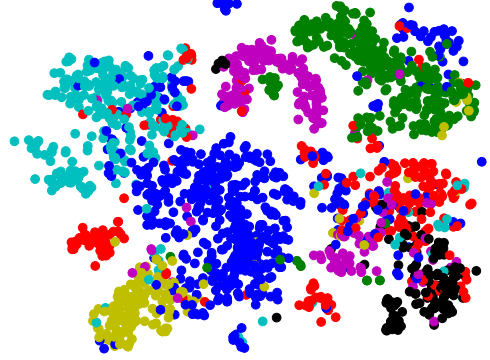
2.3 Node Embeddings

Node Embedding methods represent graph nodes in a continuous vector space. They learn a dictionary $\mathbf{Z} \in \mathbb{R}^{N \times d}$, with one d -dimensional embedding per node. Traditional methods use the adjacency matrix to learn embeddings. Skipgram models on text corpora (Mikolov et al., 2013) inspired modern graph embedding methods, which simulate random walks to learn node embeddings (Perozzi et al., 2014; Grover and Leskovec, 2016). Each random walk generates a sequence of nodes. Sequences are converted to textual paragraphs, and are passed to a word2vec-style embedding learning algorithm (Mikolov

²Self-connections added as $\mathbf{A}_{ii} = 1$, similarly to Kipf and Welling (2017).



(a) N-GCN Architecture.



(b) t-SNE visualization of fully-connected (fc) hidden layer of NGCN when trained over Cora graph.

Figure 1: (a) Model architecture, where $\hat{\mathbf{A}}$ is the normalized normalized adjacency matrix, \mathbf{I} is the identity matrix, \mathbf{X} is node features matrix, and \times is matrix-matrix multiplication. In this example, we calculate $K = 3$ powers of the $\hat{\mathbf{A}}$, feeding each power into $r = 1$ GCNs, along with \mathbf{X} . The output of all $K \times r$ GCNs can be concatenated along the column dimension, then fed into fully-connected layers, outputting C channels per node, where C is size of label space. We calculate cross entropy error between rows *prediction* $N \times C$ with known labels, and use them to update parameters of classification sub-network and all GCNs. The classification networks for N-GCN_{fc} and N-GCN_a, respectively, are described in Sections 3.3.1 and 3.3.2. (b) Demonstration of a pre-RELU activations after the first fully-connected layer of a 2-layer classification sub-network trained on Cora. Activations are PCA-ed to 50-D then visualized using t-SNE.

et al., 2013). As shown in Abu-El-Haija et al. (2018), this learning-by-simulation is equivalent, in expectation, to the decomposition of a random walk co-occurrence statistics matrix \mathcal{D} . Expectation on \mathcal{D} can be written as:

$$\mathbb{E}[\mathcal{D}] \propto \mathbb{E}_{q \sim \mathcal{Q}} [(\mathcal{T})^q] = \mathbb{E}_{q \sim \mathcal{Q}} \left[(\mathbf{D}^{-1} \mathbf{A})^q \right], \quad (3)$$

where $\mathcal{T} = \mathbf{D}^{-1} \mathbf{A}$ is the row-normalized transition matrix (a.k.a right-stochastic adjacency matrix), and \mathcal{Q} is a “context distribution” that is determined by random walk hyperparameters, such as the length of the random walk. The expectation therefore weights the importance of one node on another as a function of how well-connected they are, and the distance between them.

3 OUR METHOD

3.1 Motivation

Graph Convolutional Networks and random walk graph embeddings are individually powerful. Kipf and Welling (2017) uses GCNs for semi-supervised node classification. Instead of following traditional methods that use the graph for regularization (Belkin and Niyogi, 2003), Kipf and Welling (2017) use the adjacency matrix for training and inference, effectively diffusing information across edges at all GCN layers (see Eq. (3)). Separately, recent work has showed that random walk statistics can be

very powerful for learning an unsupervised representation of nodes that can preserve the structure of the graph (Perozzi et al., 2014; Grover and Leskovec, 2016; Abu-El-Haija et al., 2018).

Under special conditions, it is possible for the GCN model to learn random walks. In particular, consider a two-layer GCN defined in Eq. (3) with the assumption that first-layer activation is identity as $\sigma(z) = z$ and weight $\mathbf{W}^{(0)}$ is an identity matrix (either explicitly set or learned to satisfy the upstream objective). Under these two identity conditions, the model reduces to:

$$\text{GCN}_{2\text{-layer-special}}(\hat{\mathbf{A}}, \mathbf{X}) = \text{softmax} \left(\hat{\mathbf{A}}^2 \mathbf{X} \mathbf{W}^{(1)} \right) \quad (4)$$

where $\hat{\mathbf{A}}^2$ can be expanded as

$$\begin{aligned} \hat{\mathbf{A}}^2 &= \left(\mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}} \right) \left(\mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}} \right) \\ &= \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \left[\mathbf{D}^{-1} \mathbf{A} \right] \mathbf{D}^{-\frac{1}{2}} = \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathcal{T} \mathbf{D}^{-\frac{1}{2}}. \end{aligned} \quad (5)$$

By multiplying the adjacency \mathbf{A} with the transition matrix \mathcal{T} before, the $\text{GCN}_{2\text{-layer-special}}$ is effectively doing a one-step random walk, diffusing signals from nodes to neighbors without non-linearities, then applying a non-linear graph convolution layer.

3.2 Explicit Random Walks

The special conditions described above are not true in practice. Although stacking hidden GCN layers allows information to flow through graph edges, this flow is *indirect* as the information goes through feature reduction (matrix multiplication) and a non-linearity (activation function $\sigma(\cdot)$). Therefore, the vanilla GCN cannot directly learn high powers of $\hat{\mathbf{A}}$, and could struggle with modeling information across distant nodes. We hypothesize that making the GCN directly operate on random walk statistics will allow the network to better utilize information across distant nodes, in the same way that node embedding methods (e.g., DeepWalk, Perozzi et al. (2014)) operating on \mathcal{D} are superior to traditional embedding methods operating on the adjacency matrix (e.g., Eigenmaps, Belkin and Niyogi (2003)). Therefore, in addition to feeding only $\hat{\mathbf{A}}$ to the GCN model as proposed by Kipf and Welling (2017) (see Eq. (3)), we propose to feed a K -degree polynomial of $\hat{\mathbf{A}}$ to K instantiations of GCN. Generalizing Eq. (5) to arbitrary power k gives:

$$\hat{\mathbf{A}}^k = \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathcal{T}^{k-1} \mathbf{D}^{-\frac{1}{2}}. \quad (6)$$

We also define $\hat{\mathbf{A}}^0$ to be the identity matrix. Similar to Kipf and Welling (2017), we add self-connections and convert directed graphs to undirected ones, making $\hat{\mathbf{A}}$ and hence $\hat{\mathbf{A}}^k$ symmetric matrices. The eigendecomposition of symmetric matrices is real. Therefore, the low-rank approximation of the eigendecomposition Hammond et al. (2011) is still valid, and a one layer of Kipf and Welling (2017) utilizing $\hat{\mathbf{A}}^k$ should still approximate multiplication in the Fourier domain.

3.3 Network of GCNs

Consider K instantiations of $\{\text{GCN}(\hat{\mathbf{A}}^0, \mathbf{X}), \text{GCN}(\hat{\mathbf{A}}^1, \mathbf{X}), \dots, \text{GCN}(\hat{\mathbf{A}}^{K-1}, \mathbf{X})\}$. Each GCN outputs a matrix $\mathbb{R}^{N \times C_k}$, where the v -th row describes a latent representation of that particular GCN for node $v \in \mathcal{V}$, and C_k is the latent dimensionality. Though C_k can be different for each GCN, we set all C_k to be the same for simplicity. We then combine the outputs of all K GCNs and feed them into a classification sub-network, allowing us to jointly train all GCNs and the classification sub-network via backpropagation. This should allow the classification sub-network to choose features from the various GCNs, effectively allowing the overall model to learn a combination of features using the raw (normalized) adjacency, different steps of random walks (i.e. graph scales), and the input features \mathbf{X} (as they are multiplied by identity $\hat{\mathbf{A}}^0$).

3.3.1 Fully-Connected Classification Network

From a deep learning prospective, it is intuitive to represent the classification network as a fully-connected layer. We can concatenate the output of the K GCNs along the column dimension, i.e. concatenating all $\text{GCN}(\mathbf{X}, \hat{\mathbf{A}}^k)$, each $\in \mathbb{R}^{N \times C_k}$ into matrix $\in \mathbb{R}^{N \times C_K}$ where $C_K = \sum_k C_k$. We add a fully-connected layer $f_{\text{fc}}: \mathbb{R}^{N \times C_K} \rightarrow \mathbb{R}^{N \times C}$, with trainable parameter matrix $\mathbf{W}_{\text{fc}} \in \mathbb{R}^{C_K \times C}$, written as:

$$\text{N-GCN}_{\text{fc}}(\hat{\mathbf{A}}, \mathbf{A}; \mathbf{W}_{\text{fc}}, \theta) = \text{softmax} \left(\begin{array}{c} \\ \\ \\ \end{array} \left[\text{GCN}(\hat{\mathbf{A}}^0, \mathbf{X}; \theta^{(0)}) \parallel \text{GCN}(\hat{\mathbf{A}}^1, \mathbf{X}; \theta^{(1)}) \parallel \dots \right] \mathbf{W}_{\text{fc}} \right). \quad (7)$$

The classifier parameters \mathbf{W}_{fc} are jointly trained with GCN parameters $\theta = \{\theta^{(0)}, \theta^{(1)}, \dots\}$. We use subscript **fc** on N-GCN to indicate the classification network is a fully-connected layer.

3.3.2 Attention Classification Network

We also propose a classification network based on ‘‘softmax attention’’, which learns a convex combination of the GCN instantiations. Our attention model (N-GCN_a) is parametrized by vector $\mathbf{m} \in \mathbb{R}^K$, one scalar for each GCN. It can be written as:

$$\text{N-GCN}_a(\hat{\mathbf{A}}, \mathbf{X}; \mathbf{m}, \theta) = \sum_k \mathbf{m}_k \text{GCN}(\hat{\mathbf{A}}^k, \mathbf{X}; \theta^{(k)})$$

where the vector \mathbf{m} is the output of a softmax, $\mathbf{m} = \text{softmax}(\tilde{\mathbf{m}})$, and $\tilde{\mathbf{m}}$ is a vector of weights that are updated as parameters in the model.

This softmax attention is similar to ‘‘Mixture of Experts’’ model, especially if we set the number of output channels for all GCNs equal to the number of classes, as in $C_0 = C_1 = \dots = C$. This allows us to add cross entropy loss terms on all GCN outputs in addition to the loss applied at the output N-GCN, forcing all GCN’s to be independently useful. It is possible to set the $\mathbf{m} \in \mathbb{R}^K$ parameter vector ‘‘by hand’’ using the validation split, especially for reasonable K such as $K \leq 6$. One possible choice might be setting \mathbf{m}_0 to some small value and remaining $\mathbf{m}_1, \dots, \mathbf{m}_{K-1}$ to the harmonic series $\frac{1}{k}$; another choice may be linear decay $\frac{K-k}{K-1}$. These are respectively similar to the context distributions of GloVe (Pennington et al., 2014) and word2vec (Mikolov et al., 2013; Levy et al., 2015). We note that if on average a node’s information is captured by its direct or nearby neighbors, then the output of GCNs consuming lower powers of $\hat{\mathbf{A}}$ should be weighted highly.

Algorithm 1 General Implementation: Network of Graph Models

Require: $\hat{\mathbf{A}}$ is a normalization of \mathbf{A}

```
1: function NETWORK(GRAPHMODELFN,  $\hat{\mathbf{A}}$ ,  $\mathbf{X}$ ,  $L$ ,  $r = 4$ ,  $K = 6$ , CLASSIFIERFN=FCLAYER)
2:    $\mathbf{P} \leftarrow \mathbf{I}$ 
3:   GraphModels  $\leftarrow []$ 
4:   for  $k = 1$  to  $K$  do
5:     for  $i = 1$  to  $r$  do
6:       GraphModels.append(GRAPHMODELFN( $\mathbf{P}$ ,  $\mathbf{X}$ ,  $L$ ))
7:      $\mathbf{P} \leftarrow \hat{\mathbf{A}}\mathbf{P}$ 
8:   return CLASSIFIERFN(GraphModels)
```

Algorithm 2 GCN Model (Kipf and Welling, 2017)

Require: $\hat{\mathbf{A}}$ is a normalization of \mathbf{A}

```
1: function GCNMODEL( $\hat{\mathbf{A}}$ ,  $\mathbf{X}$ ,  $L$ )
2:    $\mathbf{Z} \leftarrow \mathbf{X}$ 
3:   for  $i = 1$  to  $L$  do
4:      $\mathbf{Z} \leftarrow \sigma(\hat{\mathbf{A}}\mathbf{Z}\mathbf{W}^{(i)})$ 
5:   return  $\mathbf{Z}$ 
```

Algorithm 4 N-GCN

```
1: function NGCN( $\mathbf{A}$ ,  $\mathbf{X}$ ,  $L = 2$ )
2:    $\mathbf{D} \leftarrow \mathbf{diag}(\mathbf{A}\mathbf{1})$  ▷ Sum rows
3:    $\hat{\mathbf{A}} \leftarrow \mathbf{D}^{-1/2}\mathbf{A}\mathbf{D}^{-1/2}$ 
4:   return NETWORK(GCNMODEL,  $\hat{\mathbf{A}}$ ,  $\mathbf{X}$ ,  $L$ )
```

3.4 Training

We minimize the cross entropy between our model output and the known training labels \mathbf{Y} as:

$$\min_{\Theta} \text{diag}(\mathcal{V}_L) \left[\mathbf{Y} \circ \log \text{N-GCN}(\mathbf{X}, \hat{\mathbf{A}}; \Theta) \right], \quad (8)$$

where \circ is Hadamard product, and $\text{diag}(\mathcal{V}_L)$ denotes a diagonal matrix with entry at (i, i) set to 1 if $i \in \mathcal{V}_L$ and 0 otherwise. Θ is the set of parameters to optimize. For fully-connected classification network, $\Theta = \{\mathbf{W}_{fc}, \theta\}$. For attention classification network, we apply intermediate supervision for the N-GCN_a to attempt make all GCNs become independently useful, yielding minimization objective:

$$\min_{\mathbf{m}, \theta} \text{diag}(\mathcal{V}_L) \left[\mathbf{Y} \circ \log \text{N-GCN}_a(\hat{\mathbf{A}}, \mathbf{X}; \mathbf{m}, \theta) + \sum_k \mathbf{Y} \circ \log \text{GCN}(\hat{\mathbf{A}}^k, \mathbf{X}; \theta^{(k)}) \right].$$

Algorithm 3 SAGE Model (Hamilton et al., 2017)

Require: $\hat{\mathbf{A}}$ is a normalization of \mathbf{A}

```
1: function SAGEMODEL( $\hat{\mathbf{A}}$ ,  $\mathbf{X}$ ,  $L$ )
2:    $\mathbf{Z} \leftarrow \mathbf{X}$ 
3:   for  $i = 1$  to  $L$  do
4:      $\mathbf{Z} \leftarrow \sigma(\left[ \begin{array}{c} \mathbf{Z} \\ \hat{\mathbf{A}}\mathbf{Z} \end{array} \right] \mathbf{W}^{(i)})$ 
5:      $\mathbf{Z} \leftarrow \text{L2NORMALIZEROWS}(\mathbf{Z})$ 
6:   return  $\mathbf{Z}$ 
```

Algorithm 5 N-SAGE

```
1: function NSAGE( $\mathbf{A}$ ,  $\mathbf{X}$ )
2:    $\mathbf{D} \leftarrow \mathbf{diag}(\mathbf{A}\mathbf{1})$  ▷ Sum rows
3:    $\hat{\mathbf{A}} \leftarrow \mathbf{D}^{-1}\mathbf{A}$ 
4:   return NETWORK(SAGEMODEL,  $\hat{\mathbf{A}}$ ,  $\mathbf{X}$ , 2)
```

3.5 GCN Replication Factor r

To simplify notation, our N-GCN derivations (e.g., Eq. (7)) assume that there is one GCN per $\hat{\mathbf{A}}$ power. However, our implementation feeds every $\hat{\mathbf{A}}$ to r GCN modules, as shown in Fig. 1.

3.6 Relation to other Graph Models

In addition to vanilla GCNs (e.g., Kipf and Welling, 2017), our derivation also applies to other graph models including GraphSAGE (Hamilton et al., 2017). Algorithm 1 shows a generalization that allows us to make a network of arbitrary graph models (e.g., GCN, SAGE, or others). Algorithms 2 and 3, respectively, show pseudocode for the vanilla GCN (Kipf and Welling, 2017) and GraphSAGE³ (Hamilton et al., 2017). Finally, Algorithm 4 defines our full Network of GCN model (N-GCN) by plugging Algorithm 2 into Algorithm 1. Similarly, Algo-

³Our implementation assumes mean-pool aggregation by Hamilton et al. (2017), which performs on-par to their top performer max-pool aggregation. In addition, our Algorithm 3 lists a full-batch implementation whereas (Hamilton et al., 2017) offer a mini-batch implementation.

Dataset	Type	Nodes	Edges	Classes	Features	Labeled nodes
		$ \mathcal{V} $	$ \mathcal{E} $	C	F	$ \mathcal{V}_L $
Citeseer	citaction	3,327	4,732	6 (single class)	3,703	120
Cora	citaction	2,708	5,429	7 (single class)	1,433	140
Pubmed	citaction	19,717	44,338	3 (single class)	500	60
PPI	biological	56,944	818,716	121 (multi-class)	50	44,906

Table 1: Datasets for experiments. For citation datasets, 20 training nodes per class are observed ($|\mathcal{V}_L| = 20 \times C$).

Method (Transductive)	Citeseer	Cora	Pubmed	Method (Inductive)	PPI
ManiReg (Belkin et al., 2006)	60.1	59.5	70.7	SAGE-LSTM (Hamilton et al., 2017)	61.2
SemiEmb (Weston et al., 2012)	59.6	59.0	71.1	SAGE (Hamilton et al., 2017)	60.0
LP (Zhu et al., 2003)	45.3	68.0	63.0	DCNN (our implementation)	44.0
DeepWalk (Perozzi et al., 2014)	43.2	67.2	65.3	GCN (our implementation)	46.2
ICA (Lu and Getoor, 2003)	69.1	75.1	73.9	SAGE (our implementation)	59.8
Planetoid (Yang et al., 2016)	64.7	75.7	77.2	N-GCN (ours)	46.8
GCN (Kipf and Welling, 2017)	70.3	81.5	79.0	N-SAGE (ours)	65.0
DCNN (our implementation)	71.1	81.3	79.3		
GCN (our implementation)	71.2	81.0	78.8		
SAGE (our implementation)	63.5	77.4	77.6		
N-GCN (ours)	72.2	83.0	79.5		
N-SAGE (ours)	71.0	81.8	79.4		

Table 2: Node classification performance (% accuracy for the first three citation datasets and F1 micro-averaged for multiclass PPI), using data splits of Yang et al. (2016); Kipf and Welling (2017) and Hamilton et al. (2017). We report the test accuracy corresponding to the run with the highest validation accuracy. Results above the horizontal line are copied from Kipf and Welling (2017) and from Hamilton et al. (2017) for the transductive and inductive case respectively. Because our code can recover other algorithms (as explained in Section 3.6) we show our implementations of these baselines. Finally, our proposed models are at the end of each table.

Algorithm 5 defines our N-SAGE model by plugging Algorithm 3 in Algorithm 1.

We can recover the original algorithms GCN (Kipf and Welling, 2017) and SAGE (Hamilton et al., 2017), respectively, by using Algorithms 4 (N-GCN) and 5 (N-SAGE) with $r = 1$, $K = 1$, identity CLASSIFIERFN, and modifying line 2 in Algorithm 1 to $\mathbf{P} \leftarrow \hat{\mathbf{A}}$. Moreover, we can recover original DCNN (Atwood and Towsley, 2016) by calling Algorithm 4 with $L = 1$, $r = 1$, modifying line 3 to $\hat{\mathbf{A}} \leftarrow \mathbf{D}^{-1}\mathbf{A}$, and keeping $K > 1$ as their proposed model operates on the power series of the transition matrix i.e. *unmodified* random walks, like ours.

4 EXPERIMENTS

4.1 Datasets

We experiment on three citation graph datasets: Pubmed, Citeseer, Cora, and a biological graph: Protein-Protein Interactions (PPI). We choose the aforementioned datasets because they are available online and are used by our baselines. The citation datasets are prepared

by Yang et al. (2016), and the PPI dataset is prepared by Hamilton et al. (2017). Table 1 summarizes dataset statistics.

Each node in the citation datasets represents an article published in the corresponding journal. An edge between two nodes represents a citation from one article to another, and a label represents the subject of the article. Each dataset contains a binary Bag-of-Words (BoW) feature vector for each node. The BoW are extracted from the article abstract. Therefore, the task is to predict the subject of articles, given the BoW of their abstract and the citations to other (possibly labeled) articles. Following Yang et al. (2016) and Kipf and Welling (2017), we use 20 nodes per class for training, 500 (overall) nodes for validation, and 1000 nodes for evaluation. We note that the validation set is larger than training $|\mathcal{V}_L|$ for these datasets.

The PPI graph, as processed and described by Hamilton et al. (2017), consists of 24 disjoint subgraphs, each corresponding to a different human tissue. 20 of those subgraphs are used for training, 2 for validation, and 2 for testing, as partitioned by Hamilton et al. (2017).

Method	Node per Class			
	5	10	20	100
DCNN	63.0 ± 1.0	72.3 ± 0.4	79.2 ± 0.2	82.6 ± 0.3
GCN	64.6 ± 0.3	70.0 ± 3.7	79.1 ± 0.3	81.8 ± 0.3
SAGE	69.0 ± 1.4	72.0 ± 1.3	77.2 ± 0.5	80.7 ± 0.7
N-GCN _a	65.1 ± 0.7	71.2 ± 1.1	79.7 ± 0.3	83.0 ± 0.4
N-GCN _{fc}	65.0 ± 2.1	71.7 ± 0.7	79.7 ± 0.4	82.9 ± 0.3
N-SAGE _a	66.9 ± 0.4	73.4 ± 0.7	79.0 ± 0.3	82.5 ± 0.2
N-SAGE _{fc}	70.7 ± 0.4	74.1 ± 0.8	78.5 ± 1.0	81.8 ± 0.3

Table 3: Node classification accuracy (in %) with different number of labeled nodes per class $\frac{|V|}{C} \in \{5, 10, 20, 100\}$.

4.2 Baseline Methods

For the citation datasets, we copy baseline numbers from Kipf and Welling (2017). These include label propagation (LP, Zhu et al. (2003)); semi-supervised embedding (SemiEmb, Weston et al. (2012)); manifold regularization (ManiReg, Belkin et al. (2006)); skip-gram graph embeddings (DeepWalk Perozzi et al., 2014); Iterative Classification Algorithm (ICA, Lu and Getoor, 2003); Planetoid (Yang et al., 2016); vanilla GCN (Kipf and Welling, 2017). For PPI, we copy baseline numbers from Hamilton et al. (2017), which include GraphSAGE with LSTM aggregation (SAGE-LSTM) and GraphSAGE with pooling aggregation (SAGE). Further, for all datasets, we use our implementation to run baselines DCNN (Atwood and Towsley, 2016), GCN (Kipf and Welling, 2017), and SAGE (with pooling aggregation, Hamilton et al., 2017), as these baselines can be recovered as special cases of our algorithm, as explained in Section 3.6.

4.3 Implementation

We use TensorFlow (Abadi et al., 2016) to implement our methods, which we use to also measure the performance of baselines GCN, SAGE, and DCNN. For our methods and baselines, all GCN and SAGE modules that we train are 2 layers⁴, where the first outputs 16 dimensions per node and the second outputs the number of classes (dataset-dependent). DCNN baseline has one layer and outputs 16 dimensions per node, and its channels (one per transition matrix power) are concatenated into a fully-connected layer that outputs the number of classes. We use 50% dropout and L_2 regularization of 10^{-5} for all of the aforementioned models.

⁴except as clearly indicated in Table 5

4.4 Node Classification Accuracy

Table 2 shows node classification accuracy results. We run 20 different random initializations for every model (baselines and ours), train using Adam optimizer (Ba and Kingma, 2015) with learning rate of 0.01 for 600 steps, capturing the model parameters at peak validation accuracy to avoid overfitting. For our models, we sweep our hyperparameters r , K , and choice of classification sub-network $\in \{fc, a\}$. For baselines and our models, we choose the model with the highest accuracy on validation set, and use it to record metrics on the test set in Table 2.

Table 2 shows that N-GCN outperforms GCN (Kipf and Welling, 2017) and N-SAGE improves on SAGE for all datasets, showing that *unmodified* random walks indeed help in semi-supervised node classification. Finally, our proposed models achieve state-of-the-art on all datasets.

4.5 Analysis on Graph Sparsity

We run evaluations with various sparsity level of the graph by taking different number of labeled nodes per class. In Table 3, we report node classification accuracy (in %) with various number of labeled nodes per class, $\frac{|V|}{C} \in \{5, 10, 20, 100\}$, of our models as well as several baselines. We use our largest dataset (Pubmed) for this experiment, and report mean and standard deviations on 10 runs. We use a different random seed for every run (i.e., sampling different labeled nodes), but the same 10 random seeds across models.

We see that convolution-based methods (e.g. SAGE) tend to work well with fewer training examples, while *unmodified* random walk methods (e.g. DCNN) work well with more training data. Our methods combine convolution and random walks, making them work well in both conditions.

4.6 Sensitivity Analysis

We analyze the impact of random walk length K and replication factor r on classification accuracy in Table 4. In general, model performance improves when increasing K and r . We note utilizing random walks by setting $K > 1$ improves model accuracy due to the additional information, not due to increased model capacity: Contrast $K = 1, r > 1$ (i.e. mixture of GCNs, no random walks) with $K > 1, r = 1$ (i.e. N-GCN on random walks) – in both scenarios, the model has more capacity, but the latter shows better performance. The same conclusion holds for SAGE.

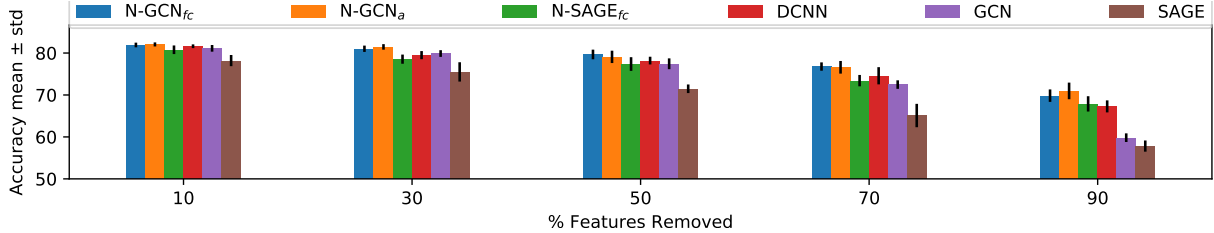


Figure 2: Classification accuracy for the Cora dataset with 20 labeled nodes per class ($|\mathcal{V}| = 20 \times C$), but features removed at random, averaging 10 runs. We use a different random seed for every run (i.e. removing different features per node), but the same 10 random seeds across models.

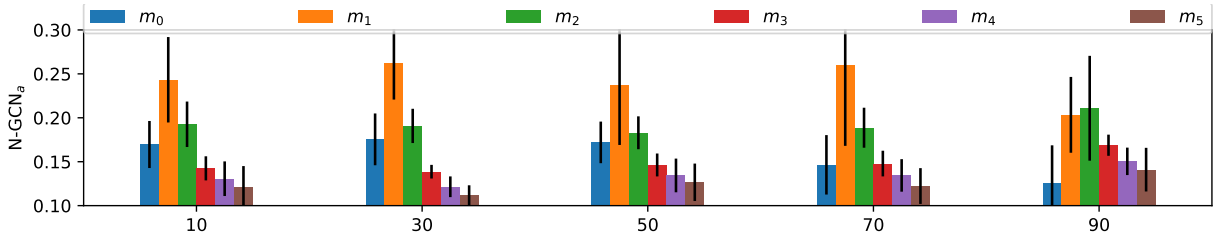


Figure 3: Attention weights (m) for N-GCN_a when trained with feature removal perturbation on the Cora dataset. Removing features shifts the attention weights to the right, suggesting the model is relying more on long range dependencies.

4.7 Tolerance to Feature Noise

We test our method under feature noise perturbations by removing node features at random. This is practical, as article authors might forget including relevant terms in the article abstract, and more generally not all nodes will have the same amount of detailed information. Fig. 2 shows that when features are removed, methods utilizing unmodified random walks (N-GCN, N-SAGE, and DCNN) outperform convolutional methods including GCN and SAGE. Moreover, the performance gap widens as we remove more features. This suggests that our methods can somewhat recover removed features by *directly* pulling-in features from nearby and distant neighbors. We visualize in Fig. 3 the attention weights as a function of features removed. With little feature removal, there is some weight on $\hat{\mathbf{A}}^0$, and the attention weights for $\hat{\mathbf{A}}^1, \hat{\mathbf{A}}^2, \dots$ decay. Maliciously dropping features causes our model to shift its attention weights towards higher powers of $\hat{\mathbf{A}}$.

4.8 Random Walk Steps vs. GCN Depth

K -step random walk will allow every node to accumulate information from its neighbors, up to distance K . Similarly, a K -layer GCN (Kipf and Welling, 2017) will do the same. The difference between the two was mathematically explained in Section 3.1. To summarize, the

former averages node feature vectors according to the random walk co-visit statistics, whereas the latter creates non-linearities and matrix multiplies at every step. So far, we display experiments where our models (N-GCN and N-SAGE) are able to use information from distant nodes (e.g. $K = 5$), but for all GCN and SAGE modules, we use 2 GCN layer for baselines and our models.

Even though the authors of GCN (Kipf and Welling, 2017) and SAGE (Hamilton et al., 2017) suggest using two GCN layers, according by holdout validation, for a fair comparison with our models, we run experiments utilizing deeper GCN and SAGE are models so that its “receptive field” is comparable to ours.

Table 5 shows test accuracies when training deeper GCN and SAGE models, using our implementation. We notice that, unlike our method which benefits from a wider “receptive field”, there is no direct correspondence between depth and improved performance.

5 RELATED WORK

We divide the rapidly-growing field of graph learning into two kinds of algorithms. The first kind is *unsupervised*, and the second kind is *supervised*.

We refer to the first kind as *embedding learning algorithms*. They input a graph and output one embed-

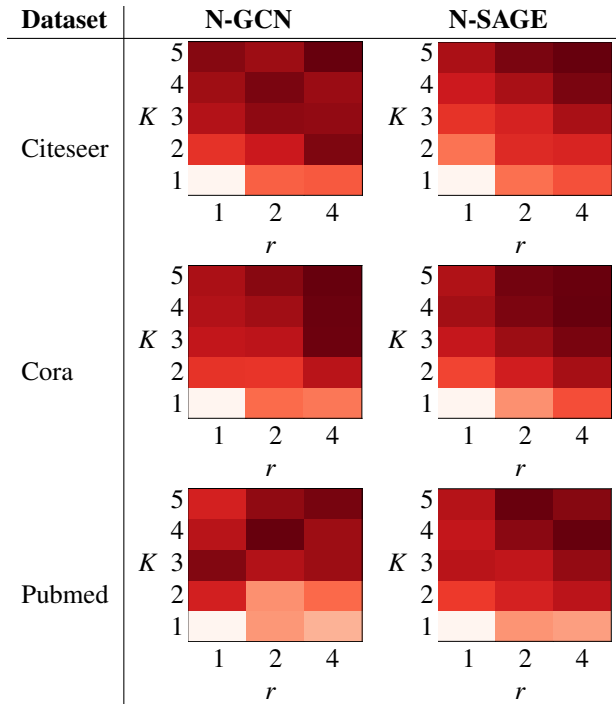


Table 4: Sensitivity Analysis: Color-coded model performance with varying random walk steps $K = \{1, 2, 3, 4, 5\}$ and replication factor $r = \{1, 2, 4\}$. Darker color means better performance, and color-scale is comparable only within the same matrix. Overall, model performance increases with larger K and r . In addition, having random walk steps (larger K) boosts performance more than increasing model capacity (larger r).

Dataset	Model	64×64		
		$64 \times C$	$\times C$	$\times 64 \times C$
Citeseer	GCN	0.699	0.632	0.659
Citeseer	SAGE	0.668	0.660	0.674
Cora	GCN	0.803	0.800	0.780
Cora	SAGE	0.761	0.763	0.757
Pubmed	GCN	0.762	0.771	0.781
Pubmed	SAGE	0.770	0.776	0.775
PPI	GCN	0.460	0.461	0.466
PPI	SAGE	0.658	0.672	0.650

Table 5: Performance of deeper GCN and SAGE models, both using our implementation. Deeper GCN (or SAGE) does not consistently improve classification accuracy, suggesting that N-GCN and N-SAGE are more performant and are easier to train. They use shallower convolution models that operate on multiple scales of the graph.

ding vector per node, with an objective of recovering the graph structure e.g. predicting every node’s neighborhood (Perozzi et al., 2014). Like ours, some of these embedding algorithms also operate on multiple scales e.g. (Cao et al., 2015; Perozzi et al., 2017; Abu-El-Haija et al., 2018). Refer to (Chen et al., 2018) for a review.

Nonetheless, our work falls under the second kind of algorithms, which directly learn a target task in a single-shot. We review methods that are most related to ours. Specifically, ones based on Graph Convolution. Defferrard et al. (2016) defines graph convolutions as a K -degree polynomial of the Laplacian, where the polynomial coefficients are learned. We differ in the order of random walk versus non-linearity. In particular, their model learns a linear combination of K -degree polynomial and pass through classifier function g , as in $g(\sum_k q_k \tilde{\mathbf{A}}^k)$, while our (e.g. N-GCN) model calculates $\sum_k q_k g(\tilde{\mathbf{A}}^k)$, where $\tilde{\mathbf{A}}$ is $\hat{\mathbf{A}}$ in our model and $\mathbf{I} - \hat{\mathbf{A}}$ in theirs, and our g can be a GCN module e.g. of (Kipf and Welling, 2017). Atwood and Towsley (2016) proposes DCNN, which calculates powers of the transition matrix and keeps each power in a separate channel until the classification sub-network at the end. Their model is therefore similar to our work in that it also falls under $\sum_k q_k g(\tilde{\mathbf{A}}^k)$. Specifically, it is a special-case of ours, when GCN module is restricted to a single layer, as explained in Section 3.6. Finally, we published a variant of this work (Abu-El-Haija et al., 2019), which inter-mixes the scales (i.e. adjacency-powers) at every graph convolutional layer: it enjoys theoretical guarantees, but its empirical performance is slightly below this version. Our open-source code implements both: <https://github.com/samihaija/mixhop>

6 CONCLUSION

We proposed a meta-model that can run arbitrary Graph Convolution models, such as GCN (Kipf and Welling, 2017) and SAGE (Hamilton et al., 2017), on the output of random walks. Traditional Graph Convolution models operate on the normalized adjacency matrix. We make multiple instantiations of such models, feeding each different graph scale i.e. a different power of the adjacency matrix, then feed the output of all instances into a classification sub-network. Our model, Network of GCNs (and similarly, Network of SAGE), is end-to-end trainable, and is able to directly learn information across near or distant neighbors. Inspecting the distribution of parameter weights in our classification sub-network, reveals that our model effectively circumvents adversarial perturbations on the input by shifting weights towards model instances operating on courser graph scales.

References

- Abadi, M. et al.
2016. TensorFlow: Large-scale machine learning on heterogeneous systems. In *Proc. of the USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, Pp. 265–283.
- Abu-El-Haija, S., B. Perozzi, R. Al-Rfou, and A. Alemi
2018. Watch your step: Learning graph embeddings through attention. In *Proc. of the Advances in Neural Information Processing Systems (NIPS)*.
- Abu-El-Haija, S., B. Perozzi, A. Kapoor, H. Harutyunyan, N. Alipourfard, K. Lerman, G. V. Steeg, and A. Galstyan
2019. Mixhop: Higher-order graph convolution architectures via sparsified neighborhood mixing. In *International Conference on Machine Learning*.
- Atwood, J. and D. Towsley
2016. Diffusion-convolutional neural networks. In *Proc. of the Advances in Neural Information Processing Systems (NIPS)*.
- Ba, J. and D. Kingma
2015. Adam: A method for stochastic optimization. In *Proc. of the International Conference on Learning Representations (ICLR)*.
- Belkin, M. and P. Niyogi
2003. Laplacian eigenmaps for dimensionality reduction and data representation. *Neural computation*, 15(6):1373–1396.
- Belkin, M., P. Niyogi, and V. Sindhwani
2006. Manifold regularization: A geometric framework for learning from labeled and unlabeled examples. *Journal of machine learning research*, 7(11):2399–2434.
- Bruna, J., W. Zaremba, A. Szlam, and Y. LeCun
2014. Spectral networks and locally connected networks on graphs. In *Proc. of the International Conference on Learning Representations (ICLR)*.
- Cao, S., W. Lu, and Q. Xu
2015. Grarep: Learning graph representations with global structural information. In *International Conference on Information and Knowledge Management*.
- Chen, H., B. Perozzi, R. Al-Rfou, and S. Skiena
2018. A tutorial on network embeddings.
- Defferrard, M., X. Bresson, and P. Vandergheynst
2016. Convolutional neural networks on graphs with fast localized spectral filtering. In *Proc. of the Advances in Neural Information Processing Systems (NIPS)*.
- Grover, A. and J. Leskovec
2016. node2vec: Scalable feature learning for networks. In *Proc. of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*.
- Hamilton, W., R. Ying, and J. Leskovec
2017. Inductive representation learning on large graphs. In *Proc. of the Advances in Neural Information Processing Systems (NIPS)*.
- Hammond, D. K., P. Vandergheynst, and R. Gribonval
2011. Wavelets on graphs via spectral graph theory. *Applied and Computational Harmonic Analysis*, 30(2):129–150.
- He, K., X. Zhang, S. Ren, and J. Sun
2016. Deep residual learning for image recognition. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Kipf, T. and M. Welling
2017. Semi-supervised classification with graph convolutional networks. In *Proc. of the International Conference on Learning Representations (ICLR)*.
- Krizhevsky, A., I. Sutskever, and G. E. Hinton
2012. ImageNet classification with deep convolutional neural networks. In *Proc. of the Advances in Neural Information Processing Systems (NIPS)*.
- LeCun, Y., L. Bottou, Y. Bengio, P. Haffner, et al.
1998. Gradient-based learning applied to document recognition. *Proc. of the IEEE*, 86(11):2278–2324.
- Levy, O., Y. Goldberg, and I. Dagan
2015. Improving distributional similarity with lessons learned from word embeddings. *Transactions of the Association for Computational Linguistics*, 3:211–225.
- Lu, Q. and L. Getoor
2003. Link-based classification. In *Proc. of the International Conference on Machine Learning (ICML)*.
- Merdan, S., C. L. Barnett, and B. T. Denton
2017. Data analytics for optimal detection of metastatic prostate cancer. Technical report, University of Michigan.
- Mikolov, T., I. Sutskever, K. Chen, G. Corrado, and J. Dean
2013. Distributed representations of words and phrases and their compositionality. In *Proc. of the Advances in Neural Information Processing Systems (NIPS)*.
- Pennington, J., R. Socher, and C. D. Manning
2014. Glove: Global vectors for word representation. In *Proc. of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Perozzi, B., R. Al-Rfou, and S. Skiena
2014. DeepWalk: Online learning of social representations. In *Proc. of the ACM SIGKDD International*

Conference on Knowledge Discovery and Data Mining.

Perozzi, B., V. Kulkarni, H. Chen, and S. Skiena

2017. Don't walk, skip!: Online learning of multi-scale network embeddings. In *International Conference on Advances in Social Networks Analysis and Mining*.

Szegedy, C., W. Liu, Y. Jia, P. Sermanet, S. E. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich

2015. Going deeper with convolutions. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

Weston, J., F. Ratle, H. Mobahi, and R. Collobert

2012. Deep learning via semi-supervised embedding. In *Neural Networks: Tricks of the Trade*, Pp. 639–655. Springer.

Yang, Z., W. Cohen, and R. Salakhutdinov

2016. Revisiting semi-supervised learning with graph embeddings. In *Proc. of the International Conference on Machine Learning (ICML)*.

Zhu, X., Z. Ghahramani, and J. Lafferty

2003. Semi-supervised learning using gaussian fields and harmonic functions. In *Proc. of the International Conference on Machine Learning (ICML)*.