

UNIVERSITY OF  
CAMBRIDGE

MATHEMATICS TRIPOS

Part III Essay

**Walking Deeper on  
Dynamic Graphs**

March 30, 2020

*Written by*  
JOSHUA SNYDER

## **Contents**

# 1 Introduction

The motivation of this essay and the methods it presents comes from the desire to understand networks of people and the relationships between them. From a mathematical perspective, the natural way to go about studying social networks is to let nodes of a graph  $G$  represent the people and edges their relationships. From this, we can induce some understanding of the neighbourhoods that a network consists of. This task of understanding the communities that exist in a social network is a complex one; even with a relatively small number of people it is not a task for which humans perform very well and the task rapidly becomes difficult as the size of the network increases.

Network embeddings are a useful way to characterise the complex structures that can exist in these graphs. The idea behind them is to represent each member of a network with a point in space, where the distance according to some metric (often simply Euclidean) between any two points corresponds to the similarity of the members in the network. These points in space are called vertex representations. The difficulty with this approach is defining a suitable similarity function  $f(p_1, p_2)$  that gives a good understanding of the network structure. This ambiguity of similarity has led to a diverse range of algorithms for calculating network embeddings.

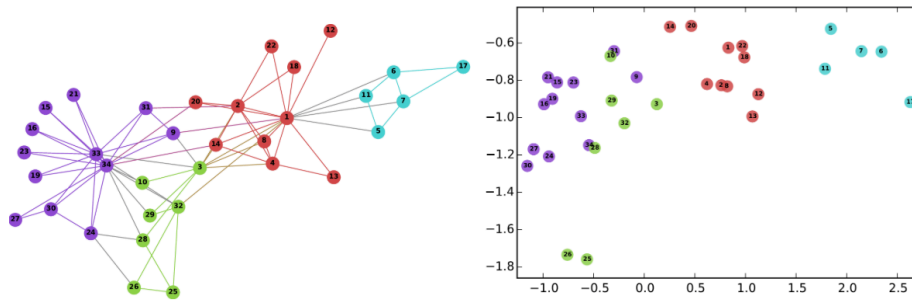


Figure 1: Two dimensional network embeddings calculated using DeepWalk[?] on Zachary’s Karate network[?]. There is a strong correspondence between community structure in the original graph and in the embedding space. Vertex colours are representing a modularity-based clustering of the input graph. In practice, network embeddings have dimensions in the tens to low-hundreds and t-SNE[?] is often used to visualise these embeddings in 2D space.

The problem of analysing social structures is made more difficult still by the fact that such networks often change over time. As new connections are formed and others die away, the network evolves and over a long time period, it can change significantly. Calculating the vertex representations for a large social graph, where the size of the network is often in the millions or billions, takes a lot of time. It is therefore computationally infeasible to recalculate these representations at every time step and so a more intelligent approach to this problem is required.

This essay looks at DeepWalk[?] and how it can be modified for dealing with dynamic graphs. DeepWalk is an algorithm that calculates vertex representations for large social graphs by using techniques adapted from natural language

processing. When proposed in 2014, the algorithm was a first of its kind and has since sparked a large amount of future research, with new algorithms improving upon it using similar methods. The core idea behind DeepWalk is to model random walks on a graph similar to sentences in a large corpus of text. This is helpful, since random walks store information about the similarity of nodes in a graph, but in a linear data type which makes it much easier to handle directly than the complex graph structures that are usually observed.

## 1.1 Summary

The section on DeepWalk explores the fundamental connection between DeepWalk and matrix factorisation as presented in the paper by Qiu et al.[?]. In this section, we demonstrate the conditions under which DeepWalk is factorising an appropriate matrix that stores some measure of similarity between nodes in the graph. In the section “Deep Drawbacks” we look at some of the criticism that followed the original DeepWalk paper and how these criticisms were addressed by future authors. This section is indicative of the shortcomings of DeepWalk but is by no means exhaustive.

After this, the section “Dynamic DeepWalking” focuses on applying a variant of the DeepWalk algorithm, proposed by Sajjad et al. [?], for dynamic graphs. This is especially important in the context of social graphs, which usually contain temporal information about user interactions. Learning social representations for dynamic graphs can allow for a better understanding of how communities move and change over time.

## 1.2 Preliminary Notation

I have endeavoured to use similar notation to the original DeepWalk paper to ease cross-referencing. However given that a number of other papers are touched upon throughout the essay, there are some pieces of notation that have been changed.

**Definition** ((Partially Labelled) Network Graph). Let  $G = (V, E)$  be an undirected graph (representing a network).  $V$  represents the members of the network, commonly referred to as the *nodes* or *vertices* and  $E \subset V \times V$  represents their connections, usually referred to as *edges*.

Some of the nodes and edges have labels and  $G_L = (V, E, X, Y)$  represents the partially labelled network. The node features are represented by  $X \in \mathbb{R}^{|V| \times S}$  where  $S$  is the size of the feature space for each attribute vector. The node labels are represented by  $Y \in \mathbb{R}^{|V| \times |\mathcal{Y}|}$ , where  $\mathcal{Y}$  is the set of labels.

Some of the nodes are labelled with  $y \in \mathcal{Y}$  and the end task is to predict the labels of the other nodes. This is called a *relational classification problem*. The aim of the DeepWalk algorithm is to learn embeddings  $W \in \mathbb{R}^{|V| \times d}$  where  $d$  is a small number of latent dimensions ( $d \ll |V|$ ) such that each of the rows in  $W$  corresponds to a  $d$  dimensional embedding of a node in  $V$ .

To avoid confusion, it is important to distinguish that the goal of DeepWalk is not to classify the nodes, but to generate useful embeddings. These embeddings are then treated as additional features and are combined with the features  $X$  to be inputted into a classification algorithm.

## 2 DeepWalk

This section gives an outline of the social representation learning algorithm DeepWalk, first introduced in the seminal paper DeepWalk: Online Learning of Social Representations by B. Perozzi et. al. [?]. The method proposed in this paper not only demonstrated performance improvements from previous methodologies but also motivated an entirely different approach. At the time of the paper being written, significant advancements were being made in natural language processing (NLP) and the idea of word embeddings was becoming popular through an embedding algorithm known as Word2Vec [?, ?]. DeepWalk implements this algorithm but replaces the idea of the context of a word in a sentence with the context of a node in a random walk on a graph. This is the crucial concept of DeepWalk from which the remaining details of the algorithm naturally follow.

The original paper on DeepWalk is lacking in a mathematical foundation and so in this section we will model the algorithm mathematically. It is suggested that the reader is familiar with the concepts outlined in the paper by Perozzi et. al. prior to reading this (more) mathematical exposition.

---

**Algorithm 1** DeepWalk

---

**Input:** network  $G(V, E)$

    window size  $k$

    embedding size  $d$

    walks per vertex  $\gamma$

    walk length  $l$

**Output:** matrix of vertex representations  $W \in \mathbb{R}^{|V| \times d}$

1: Initialization: Sample  $W$  from  $\mathcal{U}^{|V| \times d}$

2: **for**  $i = 0$  to  $\gamma$  **do**

3:      $\mathcal{O} = \text{Shuffle}(V)$

4:     **for each**  $v \in \mathcal{O}$  **do**

5:          $\mathcal{W}_v = \text{RandomWalk}(G, v, l)$

6:          $\text{Skipgram}(W, \mathcal{W}_v, k)$

7:     **end for**

8: **end for**

---

### 2.1 Introduction to DeepWalk

To understand DeepWalk mathematically, we first need to understand what the SkipGram model is doing, since this is the model underpinning DeepWalk. SkipGram was proposed by Mikolev et al. [?] as an efficient way to learn word embeddings. The idea behind SkipGram is to learn embeddings which are good at predicting nearby words in sentences. For a sentence  $w_1, \dots, w_N$ , the nearby words from  $w_i$ , known as context words, are defined as the set of words within distance  $k$ ,  $\{w_{i-k}, w_{i-k+1}, \dots, w_{i+k}\} \setminus w_i$ , where  $k$  is the size of the window. SkipGram minimises the following objective function

$$\mathcal{L} = - \sum_{c \in C(w)} \log \mathbb{P}(c|w)$$

where  $\mathbb{P}(c|w)$  is modelled by a softmax function

$$\mathbb{P}(u|v) = \frac{\exp(\vec{w} \cdot \vec{c})}{\sum_{c \in \mathcal{F}} \exp(\vec{w} \cdot \vec{c})}$$

where  $\mathcal{F}$  is the vocabulary and  $\vec{w}$  represents the embedding of the word  $w$  and  $\vec{c}$  the distributed representation of  $c$  when it serves as a context word.

\*\*\* Can I be clearer about what this distributed representation is? \*\*\*

Perozzi et al. took this model from NLP and applied it to graphs by observing that nodes in the graph can be thought of as words in an artificial language. Firstly, to find the context for each node, DeepWalk generates random walks which are analogous to sentences in a language. The context nodes are the set of nodes within a window of size  $k$ , from which input-context pairs  $(v, c)$  are constructed and added to a corpus  $\mathcal{D}$ , which is a multi-set. After these context nodes have been found the same optimization function as in SkipGram is used to learn embeddings which maximise the chance of predicting context nodes.

Algorithm 1 shows how DeepWalk is applied.  $\gamma$  represents the number of random walks that are started from each node and at the start of each pass the nodes are shuffled so that they are traversed in a random order.<sup>1</sup> For each node  $v \in V$  a random walk  $\mathcal{W}_v$  of length  $l$  is generated and used to update the network embeddings using the SkipGram algorithm applied to the random walk.

From the perspective of machine learning, SkipGram trains a neural network to do a “fake” task. It trains on this task, and then the corresponding weights learned are used as embeddings. The task is as follows:

Given an input vertex  $v$  somewhere in a random walk  $\mathcal{W}$ . Pick a nearby vertex at random. The task of the neural network is to predict the probability that each vertex in  $V$  will be this randomly chosen vertex. Therefore, vertices far away on the graph that correspond to unfamiliar nodes are unlikely to co-occur on the same random walk and will be assigned a low probability. Conversely, nearby and well connected vertices are likely to co-occur on a random walk with input  $v$  and thus will be assigned higher probabilities. This allows us to train a network with weights that represent the connectedness between nodes on the graph. The neural network is trained by feeding it pairs of nodes  $(v, c)$  where  $v$  represents the input node and  $c$  is a context node, which lies within distance  $k$  of the vertex  $v$ .

To formalise this, each of the nodes  $v \in V$  are represented by a one-hot encoding vector  $e_v \in \mathbb{R}^{|V|}$  allowing us to feed  $e_v$  into the neural network. When  $e_v$  is fed into the network, a single linear hidden layer with  $d$  neurons is used, where  $d$  is the desired dimension of the latent representations, which is then passed to a softmax classifier for output. The output of the network is a vector  $o \in \mathbb{R}^{|V|}$  containing the estimated probabilities that a randomly selected nearby word is that vocabulary word.

The idea behind having a linear hidden layer, which does not use an activation function, is to use the resulting weight matrix  $W \in \mathbb{R}^{|V| \times d}$  as the embedding vectors for the nodes in the graph. This is intuitive as the hidden layer acts as a bottleneck that tries to represent as much information as possible to distinguish

---

<sup>1</sup>This helps to prevent the algorithm from staying in a local minima when stochastic gradient descent is applied since on each iteration, as the vertices are shuffled, the shape of the loss surface changes.

the nodes, but is only allowed  $d$  neurons to do so. Since  $d \ll |V|$  there is a low risk of overfitting.

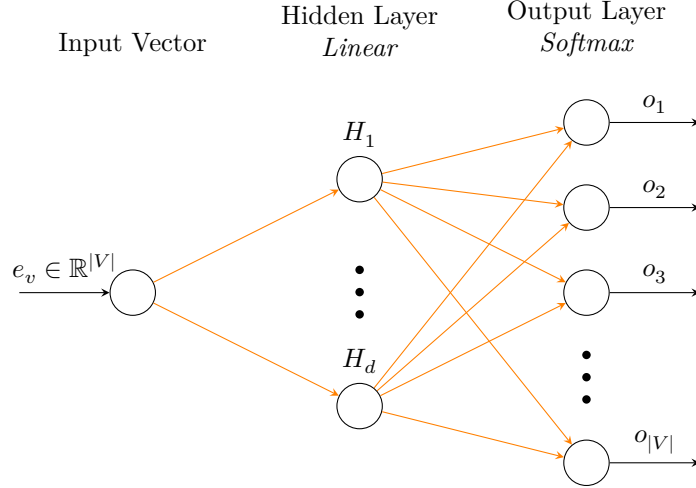


Figure 2: The structure of the neural network that is trained to obtain the weight matrix  $W$  used for the network embeddings in the SkipGram algorithm.

The algorithm used in DeepWalk varies slightly from the SkipGram algorithm discussed. Calculating the normalization factor in the Softmax layer requires a computational complexity of  $O(|V|)$ , which is intractable since this is often in the millions or billions. In the original DeepWalk paper this is reduced by using Hierarchical Softmax to approximate the softmax probabilities, requiring a complexity of only  $O(\log|V|)$ . In particular, a Huffman coding is used to reduce the access time of frequent elements in the tree, as suggested by Mikolov et al. in the original Word2Vec papers.[?, ?]

In later adaptations of DeepWalk, SkipGram with Negative Sampling (SGNS) is used instead of Hierarchical Softmax. Negative Sampling updates only a sample of the output vectors per iteration. In the remainder of this essay, when referring to DeepWalk, it will be implicitly assumed that Negative Sampling is used as appose to Hierarchical Softmax. This is because SGNS has been found to be more efficient and therefore has been adopted by much of the further literature. This convention will also serve us well when we look at applying DeepWalk to dynamic graphs since here Negative Sampling is also applied.

\*\*\* Can give a summary of NS here, if so just summarise what is said in Blog Post 2 on SkipGram \*\*\*

\*\*\* Has NS been shown, or just found empirically, to be more efficient then Hierarchical Softmax? Check node2vec paper for information on this. \*\*\*

## 2.2 SkipGram as Matrix Factorisation

In this section we will exhibit a proof that SGNS is equivalent to factorising a certain matrix  $M$  into two smaller matrices  $W$  and  $C$  where the rows in  $W$  correspond to the learned embedding of each vertex in the graph. This result

was first proved by Levy and Goldberg[?] in the context of word embeddings.

### 2.2.1 Objective of SGNS

Given an arbitrary input-context pair  $(v, c)$  the objective is to determine if the pair comes from the random-walk corpus  $\mathcal{D}$ .

Let  $P(\mathcal{D} = 1|v, c)$  denote the probability that  $(v, c)$  comes from a random walk on the graph and  $P(\mathcal{D} = 0|v, c)$  the probability it does not. Then the distribution is modelled by a sigmoid function

$$P(D = 1|v, c) = \sigma(\vec{v} \cdot \vec{c}) = \frac{1}{1 + e^{-\vec{v} \cdot \vec{c}}}$$

where  $\vec{v}$  and  $\vec{c}$  are  $d$ -dimensional vectors to be learned. SGNS attempts to maximise  $P(\mathcal{D} = 1|v, c)$  for observed pairs  $(v, c)$  whilst simultaneously maximising  $P(\mathcal{D} = 0|v, c)$  for randomly sampled negative examples.

It assumes that randomly selecting a context  $c$  for a given node  $v$  is likely to result in a “negative sample”, an unobserved pair  $(v, c)$ . In the context of social networks, this assumption is reasonable since social networks are almost always sparse (The number of edges is usually  $O(|V|)$  and not  $O(|V|^2)$ ). However in a different context, if the network is dense, then this may be an unreasonable assumption.

According to this assumption, the objective function of SGNS for a single observation  $(v, c)$  is:

$$\phi(v, c) = \log \sigma(\vec{v} \cdot \vec{c}) + b \cdot \mathbb{E}_{c_N \sim P_D} \log \sigma(-\vec{v} \cdot \vec{c})$$

where the minus sign comes from the fact that  $1 - \sigma(x) = \sigma(-x)$ ,  $b$  is the number of negative samples and  $c_N$  is the sampled context node, drawn according to  $P_D(c) = \frac{\#(c)}{|\mathcal{D}|}$  which is known as the unigram distribution.

\*\*\* The objective function and this bit generally needs more explanation, read Mikolev et al. to get a better grasp. It is not clear here. \*\*\*

**Notation.**  $\#(v, c)$ ,  $\#(v)$  and  $\#(c)$  denote the number of times vertex-context pair  $(v, c)$ , vertex  $v$  and context  $c$  appear in the generated random-walk corpus  $\mathcal{D}$  respectively.

This objective function is trained using stochastic gradient descent with updates after each observed pair in the random-walk corpus  $\mathcal{D}$ . The resulting global objective becomes

$$\mathcal{L} = \sum_{v \in V} \sum_{c \in V} \phi(v, c) \quad (1)$$

### 2.2.2 Finding the Similarity Function Learned by SkipGram

If we let  $W$  be the matrix with rows  $v_i$  (The matrix  $W$  is used to highlight that it is the weight matrix in the neural network presented in Figure 2) and  $C$  the matrix with columns  $c_i$  then SGNS can be interpreted as factorising a matrix  $M = WC^T$ . An entry in the matrix  $M_{ij}$  corresponds to the dot product  $\vec{v}_i \cdot \vec{c}_j$ . Therefore SGNS is factorising a matrix in which each row corresponds



to an input node  $v_i \in |V|$  and each column to a context node  $c_j \in |V|$  and the value of  $M_{ij}$  expresses the strength of association between the input-context pair  $(v_i, c_j)$  using some similarity function  $f(v_i, c_j)$ . What is the similarity function learned by SkipGram?

**Theorem 2.1** (Levy, Goldberg (2014)). *SkipGram with Negative Sampling (SGNS) is implicitly factorising a matrix  $M = WC^T$  with*

$$M_{ij} = \log \frac{\#(v_i, c_j) |\mathcal{D}|}{\#(v_i) \#(c_j)} - \log b$$

where  $W, C \in \mathbb{R}^{|V| \times d}$ , and  $b$  is the number of negative samples.

*Proof.* Firstly, for sufficiently large dimensionality  $d$  (so as to allow for a perfect reconstruction of  $M$ ), each of the products  $v_i \cdot c_j$  can be assumed to take their values independently of the others. \*\*\* WHY. EXPLAIN THIS. Why not sufficiently large value of  $|\mathcal{D}|$ , why does  $d$  matter here? \*\*\*

Due to this independence, the objective function  $\mathcal{L}$  can be maximised with respect to each pair  $v \cdot c$  individually.

The expectation term can be written explicitly as

$$\begin{aligned} \mathbb{E}_{\tilde{c} \sim P_{\mathcal{D}}} [\log \sigma(-\vec{v} \cdot \vec{\tilde{c}})] &= \sum_{\tilde{c} \in V} \frac{\#(\tilde{c})}{|\mathcal{D}|} \log \sigma(-\vec{v} \cdot \vec{\tilde{c}}) \\ &= \frac{\#(c)}{|\mathcal{D}|} \log \sigma(-\vec{v} \cdot \vec{c}) + \sum_{\tilde{c} \in V \setminus \{c\}} \frac{\#(\tilde{c})}{|\mathcal{D}|} \log \sigma(-\vec{v} \cdot \vec{\tilde{c}}) \end{aligned}$$

and the objective function  $\mathcal{L}$  can be expressed as

$$\mathcal{L} = \sum_{v \in V} \sum_{c \in V} \#(v, c) \log \sigma(\vec{v} \cdot \vec{c}) + \sum_{v \in V} \#(v) \left( b \cdot \mathbb{E}_{\tilde{c} \sim P_{\mathcal{D}}} [\log \sigma(-\vec{v} \cdot \vec{\tilde{c}})] \right)$$

where the second term comes from the fact that  $\#(v) = \sum_{c \in V} \#(v, c)$  by definition. Combining these equations gives that the local objective for an input-context pair is

$$\phi(v, c) = \#(v, c) \log \sigma(\vec{v} \cdot \vec{c}) + b \cdot \#(v) \cdot \frac{\#(c)}{|\mathcal{D}|} \log -\sigma(-\vec{v} \cdot \vec{c})$$

To simplify the notation let  $x = \vec{v} \cdot \vec{c}$  and, since we are assuming each of the  $\vec{v}_i \cdot \vec{c}_j$  to take their values independently, we take the partial derivative with respect to  $x$  and optimise the local objective:

$$\frac{\partial \phi}{\partial x} = \#(v, c) \cdot \log \sigma(-x) - b \cdot \#(v) \cdot \frac{\#(c)}{|\mathcal{D}|} \cdot \sigma(x)$$

Where the derivatives are since  $\sigma'(x) = \sigma(-x)$ . Setting the derivative to zero and multiplying through by  $\frac{-e^x}{\sigma(x)\sigma(-x)}$  gives:

$$\frac{b \cdot \#(v) \cdot \#(c)}{|\mathcal{D}|} e^{2x} + \left( \frac{b \cdot \#(v) \cdot \#(c)}{|\mathcal{D}|} - \#(v, c) \right) e^x - \#(v, c) = 0$$

This is a quadratic equation in  $e^x$  with two solutions. The first solution,  $e^x = -1$  is infeasible since  $x \in \mathbb{R}$  and so the appropriate solution is

$$e^x = \frac{\#(v, c) \cdot |\mathcal{D}|}{\#(v)\#(c)} \cdot \frac{1}{b}$$

Substituting  $x = \vec{v} \cdot \vec{c}$  back into the equation and taking logs gives

$$M_{ij} = \vec{v} \cdot \vec{c} = \log \left( \frac{\#(v, c) \cdot |\mathcal{D}|}{\#(v) \cdot \#(c)} \right) - \log b$$

□

Most interestingly, the resulting expression for the similarity function  $s$  is the pointwise mutual information (PMI) shifted by a factor of  $\log b$ . PMI was first introduced as a measure of association between words in 1990 by Church and Hanks [?] and became widely adopted for NLP tasks.

\*\*\* Would make sense to go into this more \*\*\*

There is an equivalent theorem for SkipGram with Softmax that was proved by Yang et al.[?] and was later used in their development of text-associated DeepWalk (TADW)[?] which incorporates text features of the vertices in a social graph.

**Theorem 2.2** (Yang et al. (2015)). *SkipGram with Softmax is implicitly factorising the matrix  $M = WC^T$  with*

$$M_{ij} = \log \frac{\#(v_i, v_j)}{\#(v_i)}$$

The proof of this follows very similarly to the previous theorem and will not be shown here since we will only be concerned with SGNS.

### 2.3 DeepWalk as Matrix Factorisation

We continue to look at DeepWalk in the context of matrix factorisation. Much of the proceeding analysis was exhibited in a recent paper by Qiu et al.[?] published in 2018 building upon work by Yang et al.[?] from 2015. The aim of the former paper was to lay the foundations for, and unify, the SkipGram based network embedding methods. First we give some preliminary definitions.

**Definition** (Adjacency Matrix (A)). The adjacency matrix  $A \in \mathbb{R}^{|V| \times |V|}$  for a graph  $G$  is the matrix with  $A_{ij} = 1$  if  $(v_i, v_j) \in E(G)$  and  $A_{ij} = 0$  otherwise.

**Definition** (Degree Matrix (D)). The degree matrix  $D \in \mathbb{R}^{|V| \times |V|}$  for a graph  $G$  is a diagonal matrix with  $D_{ii} = d_i = \text{degree}(v_i)$  for  $v_i \in V$  and  $D_{ij} = 0$  otherwise.

**Definition** (Transition Matrix (P)). The transition matrix  $P \in \mathbb{R}^{|V| \times |V|}$  for a graph  $G$  is the matrix  $P = D^{-1}A$ . It has entries  $P_{ij} = \frac{1}{d_i}$  if  $(v_i, v_j) \in E(G)$  and  $P_{ij} = 0$  otherwise. It is the transition matrix corresponding to a simple random walk on the graph  $G$ .

In their paper, Qiu et al. gave a theoretical understanding of the DeepWalk algorithm by proving the following theorem:

**Theorem 2.3** (DeepWalk as implicit matrix factorisation). *As  $l \rightarrow \infty$ , DeepWalk is equivalent to factorising*

$$\log \left( \frac{2|E|}{k} \left( \sum_{r=1}^k P^r \right) D^{-1} \right) - \log b$$

where  $b$  is the negative sampling rate.

The theorem assumes that the graph is undirected and that it is connected so that  $P$  is irreducible. It is also assumed that the graph is non-bipartite to ensure that the random walk converges to its invariant distribution. In the application to social network graphs, this assumption will almost certainly hold as the existence of an odd cycle is expected (Social networks usually contain a large amount of triangles, representing mutual friends or connections). It is possible however that the graph is not connected, in this case a dummy node can be introduced which contains edges to all nodes which will not affect the community structure, provided the graph sub-communities are sufficiently dense. What follows in the remainder of this section is a careful outline of the proof. \*\*\* If want to include details on bipartite assumption, see the folder for DeepWalk for more information. \*\*\*

Firstly,  $\pi_i = \frac{d_i}{2|E|}$  satisfies the detailed balance equations:

$$\pi_i P_{ij} = d_i \cdot \frac{1}{d_i} = d_j \cdot \frac{1}{d_j} = \pi_j P_{ji}$$

and  $\sum_{v_i \in V} \pi_i = 1$ . Thus  $\pi$  defines a distribution which is invariant for a simple random walk on the graph. Since the state space is finite and by assumption,  $P$  is irreducible, a random walk on  $G$  defines an irreducible Markov Chain  $X$  with transition matrix  $P$ . Therefore  $\pi$  is unique by the following theorem

**Theorem 2.4.** *Consider an irreducible Markov chain. Then*

- (i) *There exists an invariant distribution if and only if some state is positive recurrent.*
- (ii) *If there is an invariant distribution  $\pi$ , then every state is positive recurrent, and*

$$\pi_i = \frac{1}{\mu_i}$$

*for  $i \in S$ , where  $\mu_i$  is the mean recurrence time of  $i$ . In particular,  $\pi$  is unique.*

This is a standard proof in any course on Markov Chains and so the proof is omitted here. See Page 31 of [?] for details of a proof. As well as this, since  $P$  is irreducible and thus positive recurrent since the graph is finite, under the further assumption that a simple random walk on the graph is aperiodic the distribution of the chain tends towards  $\pi$  as  $l \rightarrow \infty$ .

To proceed with the first lemma it is useful to partition the random-walk corpus as follows:

**Definition.** For  $r = 1, \dots, k$ , we define the following

$$\mathcal{D}_{\bar{r}} = \{(v, c) : (v, c) \in \mathcal{D}, v = v_j, c = v_{j+r}\}$$

$$\mathcal{D}_{\bar{r}} = \{(v, c) : (v, c) \in \mathcal{D}, v = v_{j+r}, c = v_j\}$$

\*\*\* Need to give a proper definition to clarify what I mean here \*\*\* Thus  $\mathcal{D}_{\bar{r}}/\mathcal{D}_{\bar{r}}$  are sub-multisets of  $\mathcal{D}$  such that the context  $c$  is  $r$  steps after or before the vertex  $v$  in random walks respectively.

As an extension of previous definitions, we let  $\#(v, c)_{\bar{r}}$  and  $\#(v, c)_{\bar{r}}$  denote the number of times that an input-context pair  $(v, c)$  appears in  $\mathcal{D}_{\bar{r}}$  and  $\mathcal{D}_{\bar{r}}$  respectively. Then the following lemma holds

**Lemma 2.5.** As  $l \rightarrow \infty$ , we have

$$\frac{\#(v, c)_{\bar{r}}}{|\mathcal{D}_{\bar{r}}|} \xrightarrow{p} \pi_v(P^r)_{v,c} \text{ and } \frac{\#(v, c)_{\bar{r}}}{|\mathcal{D}_{\bar{r}}|} \xrightarrow{p} \pi_v(P^r)_{v,c}$$

*Proof.* \*\*\* A proof of this can be found in Paper 2 but the proof is not nice, perhaps I can come up with an original proof of this using convergence of the random walk to its invariant distribution.

Note that by reversibility the two statements should be provably equivalent and I have used detailed balance here to prove that the two limits given in Paper 2 are the same.

I will leave this to come back to for a nicer proof. \*\*\*

\*\*\* Note that Ref 3 gives info about as the length of the RWs becomes long, the singleton distribution of vertices will tend to the invariant distribution, referencing Modern Graph Theory, Bollobas \*\*\*  $\square$

From this we can show that

**Lemma 2.6.** As  $l \rightarrow \infty$ , we have

$$\frac{\#(v, c)}{|\mathcal{D}|} \xrightarrow{p} \frac{1}{k} \sum_{r=1}^k \pi_v(P^r)_{v,c}$$

*Proof.*

$$\begin{aligned} \frac{\#(v, c)}{|\mathcal{D}|} &= \frac{\sum_{r=1}^k (\#(v, c)_{\bar{r}} + \#(v, c)_{\bar{r}})}{\sum_{r=1}^k (|\mathcal{D}_{\bar{r}}| + |\mathcal{D}_{\bar{r}}|)} = \frac{1}{2k} \sum_{r=1}^k \left( \frac{\#(v, c)}{|\mathcal{D}_{\bar{r}}|} + \frac{\#(v, c)}{|\mathcal{D}_{\bar{r}}|} \right) \\ &\xrightarrow{p} \frac{1}{k} \sum_{r=1}^k \frac{d_v}{2|E|} (P^r)_{v,c} \end{aligned}$$

where the second equality uses the fact that  $|\mathcal{D}_{\tilde{r}}| = |\mathcal{D}_{\tilde{r}}| = \frac{|\mathcal{D}|}{2k}$  and the convergence comes from applying [Theorem XX], together with the continuous mapping theorem.  $\square$

This gives everything we need to prove the main theorem of this section, that allows us to understand the matrix that DeepWalk is implicitly factorising.

**Theorem ??** (DeepWalk as implicit matrix factorisation). *As  $l \rightarrow \infty$ , DeepWalk is equivalent to factorising*

$$\log \left( \frac{2|E|}{k} \left( \sum_{r=1}^k P^r \right) D^{-1} \right) - \log b$$

where  $b$  is the negative sampling rate.

*Proof.* Firstly, by summing the result of [Theorem XX] we get

$$\begin{aligned} \frac{\#(v)}{|\mathcal{D}|} &= \sum_{c \in V} \#(v, c) \\ &\stackrel{p}{\rightarrow} \sum_{c \in V} \frac{1}{k} \sum_{r=1}^k \pi_v(P^r)_{v,c} \\ &= \frac{\pi_v}{k} \sum_{r=1}^k \sum_{c \in V} (P^r)_{v,c} = \frac{\pi_v}{k} \sum_{r=1}^k 1 = \pi_v \end{aligned}$$

since  $P^r$  is stochastic for any  $r$ .

Similarly, using the fact that  $\pi$  is in detailed balance with  $P$ , and thus the result of [Theorem XX] can be rewritten as

$$\frac{\#(v, c)}{|\mathcal{D}|} \stackrel{p}{\rightarrow} \frac{1}{k} \sum_{r=1}^k \pi_c(P^r)_{c,v}$$

it can be shown that  $\frac{\#(c)}{|\mathcal{D}|} \stackrel{p}{\rightarrow} \pi_c$ .

Using this and by applying the continuous mapping theorem

$$\begin{aligned} \frac{\#(v, c) \cdot |\mathcal{D}|}{\#(v) \cdot \#(c)} &= \frac{\frac{\#(v, c)}{|\mathcal{D}|}}{\frac{\#(v)}{|\mathcal{D}|} \cdot \frac{\#(c)}{|\mathcal{D}|}} \stackrel{p}{\rightarrow} \frac{\frac{1}{k} \sum_{r=1}^k \frac{d_v}{2|E|} (P^r)_{v,c}}{\frac{d_v}{2|E|} \cdot \frac{d_c}{2|E|}} \\ &= \frac{2|E|}{k} \sum_{r=1}^k (P^r)_{v,c} \frac{1}{d_c} = \frac{2|E|}{k} \left( \sum_{r=1}^k (P^r D^{-1})_{v,c} \right) \end{aligned}$$

where the last equality follows since  $D$  is diagonal.

From this, applying [Theorem XX] gives that as  $l \rightarrow \infty$  DeepWalk is equivalent to factorising

$$\log \left( \frac{2|E|}{k} \left( \sum_{r=1}^k P^r \right) D^{-1} \right) - \log b$$

$\square$

This proof gives a better understanding of what the DeepWalk algorithm is doing. There are many other good questions to ask about DeepWalk, for example whether the algorithm is guaranteed to converge under stochastic gradient descent, but these questions will not be explored in this essay.

\*\*\* Here can talk about how long it takes to tend towards this. Important point, how large does  $l$  need to be for this to be meaningful? \*\*\*

## 3 Deep Drawbacks

Whilst the original DeepWalk paper made a huge impact on the future development of social representation learning on graphs, the algorithm itself has a number of drawbacks that have been pointed out in subsequent papers. In this section, I will detail some of the issues and criticism of the implementation of DeepWalk and some algorithms that were subsequently proposed to tackle these issues.

\*\*\* Plan: Pick 2 or 3 algorithms, state what they improved, link to the papers.

GraphSAGE paper points out invariance of deepwalk to an orthogonal transformation, discuss this. \*\*

### 3.1 Orthogonal Invariance

Perhaps the largest criticism that can be given to DeepWalk is inherited from its constructing part: SkipGram. As we have seen in Theorem 2.1, SkipGram is implicitly factorising a matrix  $M = WC^T$  and the embeddings  $W$  are used. We saw that by construction DeepWalk performs a similar factorisation and so do many algorithms that share similarities to DeepWalk. The learned embeddings are then combined with the features attached to each vertex and are often fed into a neural network for an upstream classification task. However the factorisation of  $M$  is not invariant under orthogonal transformation: Let  $Q \in \mathbb{R}^{d \times d}$  be an orthogonal matrix, then  $M = WC^T = WQQ^TC = (WQ)(CQ)^T$ . Therefore embeddings can undergo arbitrary orthogonal transformations without changing the resulting matrix  $M$ .

This poses a problem, since

As has been pointed out [Levy Goldberg 2015], this also calls into question the results in the original DeepWalk paper showing that DeepWalk is superior to prior network embedding learning algorithms, since the class of solutions is widely expanded for DeepWalk vs other paradigms due to this orthogonal invariance.

#### 3.1.1 Orthogonal Invariance in the Dynamic case to be put in Dynamic Drawbacks

Orthogonal invariance poses a particular problem in the case of dynamic network embeddings, since at each epoch embeddings are altered to reflect updates to the graph. However when the embeddings are updated, they can also undergo any orthogonal transformation from the previous set of embeddings and maintain the same loss. This is problematic as it makes embeddings very difficult to compare and to use for upstream classification tasks. This problem is especially prevalent when network embeddings are initialized randomly in the Unbiased Update algorithm. Whilst initializing embeddings at their previous values may help with this, over time or with the introduction of new vertices to the graph structure, orthogonal invariance remains a prevalent issue.

A possible solution to this involves careful thought of the way generated embeddings are used in upstream tasks. If the loss function for the upstream task is dependent on the network embeddings in an orthogonally invariant way then clearly orthogonal invariance is not a problem. However this reduces the ability to make comparisons between time epochs, where the potential of orthogonal

transformation must be carefully considered.

To give a concrete example of this, if  $v^s, v^t \in \mathbb{R}^d$  [Think about notation here!] represent the embeddings of a vertex at times  $s$  and  $t$  respectively, then if the upstream loss is a function of  $\|v^s\| - \|v^t\|$  then orthogonal invariance will not cause issues since  $\|Qv\| = \|v\|$  for  $Q \in \mathbb{R}^{d \times d}$  orthogonal. However a much more common loss might be a function of  $\|v^s - v^t\|$ . Such a loss function would not be appropriate, since  $\|Qv^s - v^t\| \neq \|v^s - v^t\|$  for an arbitrary orthogonal matrix  $Q$ .

A solution to this is to add a regularisation term to the loss function for Skip-Gram which penalises future embeddings based on previous ones through a metric that changes with orthogonal transformation. As an example, if  $W^t$  and  $W^{t+1}$  represent embeddings generated at times  $t$  and  $t+1$  respectively, then the euclidean distance  $\|W^{t+1} - W^t\|$  could be penalised, where the norm is restricted to vertices in  $\mathcal{V}_t \cap \mathcal{V}_{t+1}$ . This penalisation would reduce incentive for orthogonal transformations that change the embeddings across time periods, however it's effectiveness would be diminished as the time gap between two embeddings is increased. As well as this, it's affect would increase the change of getting stuck in a local minima for SGD [elaborate].

### 3.2 Embedding Space Dimensionality

The use of embeddings has become popular for a wide range of tasks primarily due to the ease with which we can manipulate mathematical expressions in ways that are seemingly impossible for more complex data structures such as graphs. In recent years there has been a large body of research devoted to algorithms for generating embeddings, but the selection of dimensionality takes a more arbitrary approach. Based on empirical evidence, a dimensionality in the low hundreds is often suggested[?, ?].

The performance of network embeddings is sensitive to it's dimensionality in a number of ways. Having too large a dimension increases the likelihood of overfitting, thus increasing the variance of the embeddings with respect to small updates in the graph. This can significantly worsen the ability of network embeddings to give a faithful representation of the graph structure and reduces the performance of downstream classification tasks. As well as this, since an increase in the dimensionality of network embeddings effects the number of model parameters, it increases learning time. This means that, in practice embedding dimensionality is often reduced based on computational savings rather than hyperparameter selection, which can severely reduce the embeddings ability to distinguish between graphs effectively, especially in the dynamic case where graph updates are small and incremental.

[Hammer into this point, but from a more DeepWalk specific viewpoint. The main point here is that specifically in this dynamic application, the small incrementing of graphs over time results in a high sensitivity to the embedding space dimensionality. 1803.00502v4, Yin is a great paper for looking more at this.]

### 3.3 Exploring Local and Global Properties of Graphs

There are two extreme strategies to search a graph starting from a vertex. The first is Breadth-First Search (BFS) which samples the immediate neighbours of



the node and helps to understand local structure within the nodes close proximity. The other is Depth-First Search (DFS) which samples nodes by going as far away as possible from the starting vertex. BFS focusses on understanding local structure and DFS on global structure. By using unbiased random walks to try and get a fair coverage of the graph, DeepWalk does not allow for orientating the focus to a depth-first or breadth-first approach which is useful for many applications. A different algorithm, node2vec, was proposed by Grover and Leskovec[?] and aims to combine the properties of BFS and DFS through introducing biased random walks. New parameters  $p, q$  are introduced which control the transition probabilities of the random walk on the graph. The unnormalised transition probability from  $y$  to  $z$  given that the walk transitioned to  $y$  from node  $x$  is:

$$p_x(y, z) = \begin{cases} \frac{1}{p} & \text{if } z = x \\ 1 & \text{if } d(x, z) = 1 \\ \frac{1}{q} & \text{if } d(x, z) = 2 \end{cases} \quad (2)$$

Varying  $p$  and  $q$  affects the random walks likelihood to travel far away versus exploring locally. Therefore by controlling these parameters a balance between BFS and DFS can be achieved. By biasing the random walks in a creative way and learning good parameters for  $p, q$  appropriately for each graph (e.g using a grid search) node2vec shows superior performance for learning vertex representations used in a number of upstream classification tasks, including on the BlogCatalog and Wikipedia datasets used in the original DeepWalk paper. Since DeepWalk uses unbiased random walks, it corresponds to using node2vec with  $p = q = 1$ . Whilst node2vec brings more flexibility, it still suffers from the orthogonal invariance problem. As well as this, it is not always clear if the values of  $p$  and  $q$  give a meaningful interpretation of the importance of BFS and DFS in the exploration of a given network, or whether performance enhancements are simply gained from an increase in the dimension of the hyperparameter space.

### 3.4 Capturing Structural Similarity

Another drawback of DeepWalk is that it does well at capturing neighbourhoods in the graph, but does not perform well at capturing structural similarities amongst nodes. That is to say that DeepWalk performs well at identifying node similarity via homophily but not by structural equivalence. Applying DeepWalk to two nodes that have a similar structural role in the graph but are in faraway locations will result in embeddings that are distant from each other. An alternative approach that performs significantly better than both DeepWalk and node2vec at this task was suggested by Ribeiro et al. in their paper on struc2vec[?]. Unlike DeepWalk, struc2vec does not take the distance between nodes into account and instead builds a context graph where nodes are close together if they are structurally similar and then SkipGram is applied similarly to DeepWalk, but this time on the context graph instead of the original graph. The approach of struc2vec is clearly to recognise a different kind of network structure to the one DeepWalk was created to recognise. struc2vec is more applicable to classification tasks where the labels depend heavily on the structural role of the nodes. As an example, struc2vec would be more suitable for the task of finding social group leaders within a large network graph. In such a task it is not the communities that are important to recognise, but the roles people

play within their communities. For this reason, it is important to identify the suitability of a relational classification problem to the application of DeepWalk.

## 4 Dynamic DeepWalking

In this section we transition away from the static implementation of DeepWalk towards network embeddings for Dynamic datasets. In particular, the focus of the section is to introduce an adaptation of DeepWalk to dynamic datasets as proposed in the paper published by Sajjad et al.[?] in early 2019.

As stated in the motivation for the paper, most of the recently developed representation learning methods can only be applied to static graphs while many real-world graphs are constantly changing over time. Therefore, to apply these methods the graph must be reanalysed at regular snapshots in time. This is very inefficient. Sajjad et al. proposed a modified version of the DeepWalk algorithm that is better suited to dynamic graphs and utilises the fact that the graph structure is unlikely to change significantly to develop a much more efficient way of analyzing social representation on dynamic graphs. The computational complexity of the algorithm depends on the graph density and number of edges added per epoch, however in the case of social networks these are both low and thus the algorithm is well suited to this application.

Previous algorithms to be used on dynamic graphs have not used what was learned about the graph in previous snapshots to efficiently calculate the node representations in the next snapshot of time. The problem of doing so can be split into two main tasks:

1. Generate a random walks corpus for a new snapshot of the graph based on a corpus from an older snapshot. In particular, we seek to find an efficient update algorithm that generates a set of random walks on the new snapshot from random walks on the old one, that is statistically indistinguishable from generating a set of random walks from scratch on the new snapshot of the graph.
2. Update the vertex representations incrementally so that they do not need to be generated from scratch every time step.

In what follows both of these tasks are explored. Firstly it is necessary to introduce notation appropriate for dynamic graphs.

**Definition.** A dynamic network graph is a series of network graphs  $\mathcal{G}^t = \{\mathcal{V}^t, \mathcal{E}^t\}$  with  $\mathcal{V}^t = \{v_1^t, \dots, v_{n(t)}^t\}$  and  $\mathcal{E}^t = \{e_1^t, \dots, e_{m(t)}^t\}$ , with  $t$  a discrete series of times.

A dynamic network graph is a series of updates to network graphs that consist of adding/removing vertices and/or edges to the graph.

**Definition** (Deleted and added vertex sets). The sets of vertices and edges of the graph that are deleted from time  $t$  to  $t + 1$  are denoted by  $D_{\mathcal{V}}^{t+1}$  and  $D_{\mathcal{E}}^{t+1}$  respectively. The sets of added vertices and edges are denoted  $A_{\mathcal{V}}^{t+1}$  and  $A_{\mathcal{E}}^{t+1}$ .

The sets of all vertices contained in the added and deleted edges are denoted  $\mathcal{V}(A_{\mathcal{V}}^{t+1})$  and  $\mathcal{V}(D_{\mathcal{E}}^{t+1})$  respectively.

Then to analyse the effect of changes on the random walks we define the following terms:

- *Affected vertices*: All vertices that are in an added or removed edge which have not been deleted.

$$\mathcal{V}_{\text{affected}}^{t+1} = \mathcal{V}(A_{\mathcal{E}}^{t+1}) \cup \mathcal{V}(D_{\mathcal{E}}^{t+1}) \setminus D_{\mathcal{V}}^{t+1}$$

- *Affected walks*: All random walks in the corpus of random walks  $W^t$  that contain at least one affected vertex.

Importantly all unaffected walks in the corpus  $W^t$  remain valid random walks on the graph  $\mathcal{G}_{t+1}$

### 4.1 Updating the Random Walk Corpus

As a preliminary solution to the problem of generating new random walks for the random walk corpus an algorithm called Naïve Update Algorithm 2 is introduced. This algorithm generates random walks of length  $l$  for each of the affected vertices. Once this is complete, it updates the random walks  $W^t$  by replacing the old walks with the corresponding re-generated walks and adds the random walks initiated from new vertices.

However, Sajjad et al. show by way of example that the Naïve update algo-

---

#### Algorithm 2 Naïve Update

---

```

1: procedure NAÏVEUPDATE( $\mathcal{G}^{t+1}, W^t, \mathcal{V}_{\text{affected}}^1, r, l$ )
2:   ▷ Initialise  $r$  new walks from each of the affected vertices in  $\mathcal{V}_{\text{affected}}^1$ 
3:    $W \leftarrow \text{initWalks}(\mathcal{V}_{\text{affected}}^1, r)$ 
4:   ▷ *** WHAT DOES THIS ACTUALLY DO? ***
5:    $W^{t+1} \leftarrow \text{randomwalk}(W, l, \mathcal{G}^{t+1})$ 
6:   ▷ Replace the old walks and add new ones
7:    $W^{t+1} \leftarrow \text{update}(W^t, W^{t+1})$ 
8:   return  $W^{t+1}$ 
9: end procedure

```

---

rithm has biased empirical transition probabilities and that the resulting corpus of random walks is not statistically indistinguishable from generating a new random walk corpus at time  $t + 1$ . To solve this problem a more sophisticated algorithm must be introduced.

#### 4.1.1 Unbiased Update

\*\*\* In the procedures I need to be clear about what `initWalks` and `randomwalk` are doing \*\*\*

The motivation behind the unbiased update algorithm is the following: Consider a random walk that has arrived at a vertex  $v$  at any point along the walk, if this vertex is not in the set  $\mathcal{V}_{\text{affected}}^{t+1}$  then there is no change to it's neighbours and thus the choice for the next vertex in the random walk remains the same. However if  $v \in \mathcal{V}_{\text{affected}}^{t+1}$  then the neighbours of the vertex have changed which makes the random walk biased from this vertex onwards since it evolves with incorrect transition probabilities. Therefore only the affected walks need to be updated to retain a statistically satisfactory set of random walks at time  $t + 1$ .

**Algorithm 3** Unbiased Update

---

```

1: procedure UNBIASEDUPDATE( $\mathcal{G}^{t+1}, W^t, \mathcal{V}_{\text{affected}}^1, r, l$ )
2:    $\triangleright$  Partition the affected vertices into new and existing ones
3:    $\mathcal{V}_n \leftarrow \text{newVertices}(\mathcal{V}_{\text{affected}}^1)$ 
4:    $\mathcal{V}_e \leftarrow \text{existingVertices}(\mathcal{V}_{\text{affected}}^1)$ 
5:    $\triangleright$  Filter the existing random walk corpus to get only the walks that
      contain affected vertices (in  $\mathcal{V}_e$ )
6:    $W_{\text{affected}} \leftarrow \text{filter}(W^t, \mathcal{V}_e)$ 
7:    $\triangleright$  Trim the affected walks to the first affected vertex
8:    $W_e \leftarrow \text{trim}(W_{\text{affected}}, \mathcal{V}_e)$ 
9:    $\triangleright \leftarrow$  Initialise  $r$  walks from each of the new vertices in  $\mathcal{V}_n$ 
10:   $W_n \leftarrow \text{initWalks}(\mathcal{V}_n, r)$ 
11:   $\triangleright$  Take the union of the new walks to create the updated corpus
12:   $W \leftarrow W_e \cup W_n$ 
13:   $W^{t+1} \leftarrow \text{randomwalk}(W, l, \mathcal{G}^{t+1})$ 
14:   $W^{t+1} \leftarrow \text{update}(W^t, W^{t+1})$ 
15:  return  $W^{t+1}$ 
16: end procedure

```

---

The Unbiased Update algorithm updates only those walks that contain vertices in  $\mathcal{V}_{\text{affected}}^{t+1}$  and does so by re-sampling the affected walks from the first affected vertex in the walk. Searching for these affected vertices is computationally expensive so another algorithm, Fast Update, is also suggested that generates the affected walks from the beginning instead of trimming them. Whilst computationally more efficient, this generates a biased random walk corpus. This is because by applying the Fast Update routine, we have effectively generated random walks on the graph  $\mathcal{G}^{t+1}$  and then regenerated all walks that hit an affected vertex, results in a corpus that is biased towards walks that do not visit the affected vertices.

In practice, when implementing the Unbiased Update algorithm the random seed that is used to generate each of the walks can be stored and then used to regenerate the random walks up to their first affected vertex by reusing the seed. Therefore the Unbiased Update algorithm can be implemented with the same complexity as the Fast Update algorithm.

\*\*\* Can I give a proof that the Unbiased Update algorithm does indeed generate an unbiased corpus? They have given an explanation but it is not very mathematical. \*\*\*

\*\*\* Can include the complexity analysis if need be but it is not particularly informative \*\*\*

## 4.2 Updating the Vertex Representations Efficiently

Once the updated random walk corpus  $W^{t+1}$  has been calculated, the vertex representations must be updated. In the DeepWalk algorithm, SkipGram is used to optimize the objective function  $\mathcal{L}$  using SGNS. The algorithm starts by initialising the vertex representations randomly and uses stochastic gradient descent to optimize the objective function over the input-context pairs in  $\mathcal{D}_{t+1}$ . However, initialising the representations randomly in the dynamic case is inef-

ficient since between each epoch the optimal vertex representations are similar, provided only a small portion of the graph has changed. It is therefore much more efficient to initialise the vertex representations at time  $t + 1$  as the corresponding representations at time  $t$  with any new vertices being randomly initiated, before performing stochastic gradient descent to minimize the SGNS objective function over the new random walks corpus  $W^{t+1}$ . The idea behind this is similar to using pretrained weights in a transfer learning problem instead of initialising weights randomly and results in stochastic gradient descent converging at a much faster rate.

A potential downside of this is that if the objective function is not convex and stochastic gradient descent gets stuck in a local minima whilst training then it is more likely to stay there. To reduce the likelihood of getting stuck in a local minima, at each epoch a small random subset of embedding coordinates could be chosen to be randomly initialised.

## 5 Conclusion

See the essay descriptors. It is required to have a conclusion and suggested to reccomend further areas of research!