#### Adversarial Attacks on Node Embeddings via Graph Poisoning

Aleksandar Bojchevski, Stephan Günnemann

Technical University of Munich

ICML 2019



#### $\sum_{k}$ Node embeddings are used to

- Classify scientific papers
- Recommend items
- Classify proteins
- Detect fraud
- Predict disease-gene associations
- Spam filtering



num. papers

### <u>ہہ</u> Background: Node embeddings

Every node  $v \in \mathcal{V}$  is mapped to a low-dimensional vector  $z_v \in \mathbb{R}^d$  such that the graph structure is captured.



Similar nodes are close to each other in the embedding space.

Background: Random walk based embeddings

Let nodes = words and random walks = sentences.

Train a sequence model, e.g. Word2Vec.



Nodes that co-occur in the random-walks have similar embeddings.

# Are node embeddings robust to adversarial attacks?

In domains where graph embeddings are used (e.g. the Web) adversaries are common and false data is easy to inject.

### Adversarial attacks in the graph domain



clean graph

poisoned graph



The optimal embedding from the to be optimized graph G

### Poisoning attack for random walk models



The optimal embedding from the to be optimized graph G

### $G_{pois.} = \underset{\substack{G \in all \ graphs}{|G_{clean} - G| \le budget}}{\operatorname{argmax}} \underset{\substack{Z \\ T_1, T_2, \dots, T_G, Z}{\min \mathcal{L}(\{r_1, r_2, \dots, T_G, Z\})}$

#### Challenges

Bi-level optimization problem.

Combinatorial search space. + - - - -

Inner optimization includes non-differentiable sampling.





(1a) Matrix factorization

(1b) optimal  $\mathcal L$  via spectrum

(2) Approximate poisoned spectrum

#### Overview

- 1. Reduce the bi-level problem to a single-level
  - a) DeepWalk as Matrix Factorization
  - b) Express the optimal  $\mathcal{L}$  via the graph spectrum
- 2. Approximate the poisoned graph's spectrum





## 1. Reduce bi-level problem to a single-level

a) DeepWalk corresponds to factorizing the PPMI matrix.

$$M_{ij} = \log \max\{cS_{ij}, 1\} \qquad S = (\sum_{r=1}^{T} P_{\uparrow}^{r}) D^{-1}$$

Get the embeddings Z via SVD of M

Rewrite *S* in terms of the generalized spectrum of *A*.

$$Au = \lambda Du \qquad S = U \left( \sum_{r=1}^{T} \Lambda^r \right) U^T$$

generalized eigenvalues/vectors

![](_page_12_Picture_0.jpeg)

### م 1. Reduce bi-level problem to a single-level

b) The optimal loss is now a simple function of the eigenvalues.

$$\min_{Z} \mathcal{L}(G, Z) = f(\lambda_i, \lambda_{i+1}, \dots)$$

Training the embedding is replaced by computing eigenvalues.

$$G_{pois.} = \operatorname*{argmax}_{G} \min_{Z} \mathcal{L}(G, Z) \implies G_{pois.} = \operatorname*{argmax}_{G} f(\lambda_{i}, \lambda_{i+1}, ...)$$

![](_page_13_Figure_0.jpeg)

### 2. Approximate the poisoned graph's spectrum

Compute the change using Eigenvalue Perturbation Theory.

$$\begin{aligned} A_{poisoned} &= A_{clean} + \Delta A \\ -\lambda_{poisoned} &= \lambda_{clean} + u_{clean}^{T} (\Delta A + \lambda_{clean} \Delta D) u_{clean} \\ \text{simplifies for a single edge flip } (i, j) \\ \lambda_{p} &= \lambda_{c} + \Delta A_{ij} (2u_{ci} \cdot u_{cj} - \lambda_{c} (u_{ci}^{2} + u_{cj}^{2})) \\ \end{aligned}$$

![](_page_14_Figure_0.jpeg)

![](_page_14_Figure_1.jpeg)

(1a) Matrix factorization

(1b) optimal  $\mathcal L$  via spectrum

(2) Approximate poisoned spectrum

#### Overall algorithm

- 1. Compute generalized eigenvalues/vectors ( $\Lambda/U$ ) of the graph
- 2. For all candidate edge flips (i,j) compute the change in  $\lambda_i$
- 3. Greedily pick the top candidates leading to largest optimal loss

![](_page_15_Picture_0.jpeg)

Poisoning decreases the overall quality of the embeddings.

![](_page_15_Figure_2.jpeg)

Our attacks:	
Gradient baseline:	
Simple baselines:	$\star \times \bullet$

Clean graph: ----

Targeted attack

#### Goal: attack a specific node and/or a specific downstream task.

![](_page_16_Figure_2.jpeg)

- Misclassify a single given target node t - - -
- Increase/decrease the similarity of a set of node pairs  $\mathcal{T} \subset \mathcal{V} \times \mathcal{V}$

Targeted attack

Most nodes can be misclassified with few adversarial edges.

![](_page_17_Figure_2.jpeg)

Transferability

Our selected adversarial edges transfer to other (un)supervised methods.

budget	DeepWalk SVD	DeepWalk Sampling	node 2vec	Spectral Embed.	Label Prop.	Graph Conv.
250	-7.59	-5.73	-6.45	-3.58	-4.99	-2.21
500	-9.68	-11.47	-10.24	-4.57	-6.27	-8.61

The change in  $F_1$  score (in percentage points) compared to the clean graph. Lower is better.

### Analysis of adversarial edges

There is no simple heuristic that can find the adversarial edges.

![](_page_19_Figure_2.jpeg)

Poster: #61, Pacific Ballroom, Today

Code: github.com/abojchevski/node\_embedding\_attack

![](_page_20_Figure_2.jpeg)

![](_page_20_Figure_3.jpeg)

 $\min_{Z} \mathcal{L} = f(\square)$ 

	+		$\approx$					

(1a) Matrix factorization

(1b) optimal  ${\cal L}$  via spectrum

(2) Approximate poisoned spectrum

#### Summary

□ Node embeddings are vulnerable to adversarial attacks.

- □ Find adversarial edges via matrix factorization and the graph spectrum.
- Relatively few perturbations degrade the embedding quality and the performance on downstream tasks.

![](_page_20_Picture_15.jpeg)