## Adversarial Attacks on Neural Networks for Graph Data

**Daniel Zügner**, Amir Akbarnejad, Stephan Günnemann Technical University of Munich, Germany

KDD 2018

Deep Learning on Graphs

Graphs are ubiquitous

- Social networks
- Web graphs
- Knowledge graphs



- (Unsupervised) learning of node representations (e.g. DeepWalk)
- Neural message passing (graph convolution)
- Implicit generative models for graphs



## Attacks on Deep Learning for Graphs

Adversaries are very common in application scenarios, e.g. search engines, or recommender systems.

These adversaries will exploit any vulnerabilities exposed.

In our work, we try to answer the question:

## Are deep learning models for graphs robust with respect to adversarial attacks?





Partially labeled, attributed graph

Message passing (a.k.a. graph convolution) aggregates local information. This could mean **higher robustness** – or lower robustness due to **cascading failures**!



• Image of a tabby cat correctly classified





- Image of a tabby cat correctly classified
- Add imperceptible perturbation
- Model classifies the cat as guacamole



Perturbation



### Adversarial attacks are a **real threat**. When dealing with **graphs**, we need to **rethink possible attack** and defense **strategies**.





#### Adversarial Attacks – Graphs Model Training **Original** graph Сору **Evasion attack** Adversarial attack Reuse Updated **Poisoning attack** model Modified graph **Re-train**

**Transductive learning**: data consists of labeled and unlabeled samples; all data used for training.

### **Evasion attack**: Modify data to fool a static classifier.

But: modifications are on the training data (transductive setting).Re-training can restore the predictions







# Poisoning attack on node classification

$$\arg \max_{A',X'} \max_{c \neq c_{old}} \log Z^*_{v,c} - \log Z^*_{v,c_{old}}$$

where 
$$Z^* = f_{\theta^*}(A', X') = softmax(\hat{A}' \operatorname{ReLU}(\hat{A}'X'W^{(1)})W^{(2)}),$$
  
with  $\theta^* = \arg\min_{\theta} \mathcal{L}(\theta; A', X')$ 

 $s.t.(A',X') \approx (A,X)$ 

 $A \in \{0,1\}^{N \times N}$ : original adjacency matrix  $X \in \{0,1\}^{N \times D}$ : (binary) node attributes A': modified structure X': modified features

v : target node

Poisoning attack on node classification  $\arg\max_{A',X'}\max_{c\neq c_{old}}\log Z^*_{v,c} - \log \overline{Z^*_{v,c_{old}}}$ Message passing where  $Z^* = f_{\theta^*}(A', X') = softmax(\hat{A}' ReLU(\hat{A}' X' W^{(1)}) W^{(2)}),$ with  $\theta^* = \arg \min_{\theta} \mathcal{L}(\theta; A', X')$  (after re-train) c.f.  $\mathcal{L}(\theta; A, X)$ : evasion **Unnoticeability**"  $A \in \{0,1\}^{N \times N}$ : original adjacency matrix  $s.t.(A',X') \approx (A,X)$  $X \in \{0,1\}^{N \times D}$ : (binary) node attributes constraint A': modified structure X': modified features v : target node



# Challenges in the Graph Setting

1. arg max: optimization over **discrete variables** (gradient information less reliable)

2. Relational dependencies between the nodes: propagation effects



4.  $\theta^* = \arg \min_{\theta} \mathcal{L}(\theta; A', X')$ : minimize classification accuracy **after** (re-)**training** on the modified data (transductive setting)





Our approach

- Use linear surrogate model to perform attacks efficiently while enforcing unnoticeability constraints on the changes.
- Train state-of-the-art models on the perturbed data and evaluate the degradation in performance.
- No access to the classifier is needed for the attack.



Based on a two-layer Graph Convolutional Network (GCN):

$$Z = f_{\theta}(A, X) = softmax \left( \hat{A} R \mathcal{O} U \left( \hat{A} X W^{(1)} \right) W^{(2)} \right)$$
 Linearize classifier

$$\log Z' = \hat{A}^2 X W'$$

**Structure** perturbations:

**Feature** perturbations:

$$\max_{\hat{A}} \mathcal{L}'(\log Z'_{\nu}) \text{ where } \log Z'_{\nu} = [\hat{A}^2 C]_{\nu} \text{ Constants}$$
$$\max_{X} \mathcal{L}'(\log Z'_{\nu}) \text{ where } \log Z'_{\nu} = [C_1 X C_2]_{\nu}$$





Based on a two-layer Graph Convolutional Network (GCN):

In contrast to the gradient, this simpler **surrogate model** allows us to **analytically** and **efficiently** determine the **exact impact** of a perturbation.

Structure perturbations:  $\max_{\widehat{A}} \mathcal{L}'(\log Z'_v)$  where  $\log Z'_v = [A^2 C]_v$  Constants Feature perturbations:  $\max_{X} \mathcal{L}'(\log Z'_v)$  where  $\log Z'_v = [C_1 X C_2]_v$ 



### Unnoticeability Constraint

 $(A', X') \approx (A, X)$ : Visual inspection by a human is not an option for graphs.

What are sensible measures of 'closeness' for graphs?

Structure perturbations:  $A' \approx A$ 

**Statistical test** on the original and modified **degree distributions** to ensure structural similarity.



Feature perturbations:  $X' \approx X$ 

**Co-occurrence constraint** for features to prevent addition of unrealistic, easy to detect features

Adversarially inserted words

to ML paper abstracts:

with constraint	without constraint	
probabilistic	efforts	
bayesian	david	
inference	family	

### ر Unnoticeability Constraint م

 $(A', X') \approx (A, X)$ : Visual inspection by a human is not an option for graphs. What are sensible measures of 'closeness' for graphs?

Structure perturbations:  $A' \approx A$  Feature perturbations:  $X' \approx X$ 

Statistical test on the original and modified degree Co-occur

## **Closed-form** equations, enabling a check for violations of the constraints in **constant time**.



Dataset	Nodes	Edges	Attributes
Cora-ML	2,485	5,069	1,433
Citeseer	2110	3,668	3,703
PolBlogs	1,222	16,714	_



- Experimental evaluation on three datasets: Cora (ML), Citeseer, PolBlogs
- Transfer experiments with Graph Convolutional Network (GCN), Column Network, and DeepWalk
- **Perturbation budget is d+2**, where d is the target's degree
- Evaluation on **5 different splits**; **10x re-training** per attack

Baselines:

- Inter-class random (direct; structure): insert edges randomly to nodes from different classes.
- Gradient (direct; structure): insert/remove edges based on the gradient.







### Results: Example Attack on GCN

Predicted probabilities after **attack** (5 modifications)



<u>ہ Results:</u> Transfer

#### Poisoning attack on **GCN**



#### Poisoning attack on **DeepWalk**



#### Deep learning models for graphs are **not robust** to adversarial attacks.

# Results: Limited knowledge

**Setup:** Only provide a **small part of the network** around the target node to the surrogate model to attack (evaluation with GCN on the complete graph).







- Deep learning models for graphs are highly vulnerable to adversarial attacks.
- We propose an efficient algorithm for performing transferable attacks.
- These attacks are **successful** even under **restrictive attack scenarios**, e.g. no access to target node or limited knowledge about the graph.





github.com/danielzuegner/nettack

