PROVABLY ROBUST MACHINE LEARNING ON GRAPHS

Aleksandar Bojchevski 22.03.202 I NEC Labs Europe

CAN YOU TRUST A HORSE?

1-1 or 12 or 13 ou 14 00 15 an 15 &

32 1 3.3 1 3.4 1 35.

w 4.2043 04.40th

MAKE SURE YOUR MODEL IS NOT A HORSE!

GRAPHS ARE EVERYWHERE



THE GRAPH ABSTRACTION



GRAPH-BASED MODELS ARE USED TO



COMMON ASSUMPTIONS



observed graph

NOISE



The graph is an output of a complex pipeline

NOISE AND ADVERSARIES





observed <mark>perturbed</mark> graph

latent clean graph

NOISE AND ADVERSARIES





latent clean graph

observed perturbed graph

Machine learning on graphs in real-world settings

FEATURES CAN BE PERTURBED TOO





PROBLEM SETUP

Semi-supervised node classification



PROBLEM SETUP

Semi-supervised node classification





The prediction is (?)

THE ADVERSARY WANTS TO

Change the prediction of a target node



The prediction is (1) after perturbation

HOW EASY IS TO MANIPULATE THE PREDICTION?



Attackers can misclassify most nodes by perturbing just a few edges

Zügner, Akbarnejad, Günnemann. "Adversarial Attacks on Neural Networks for Graph Data". KDD 2018.

ROBUSTNESS CERTIFICATE

Provable guarantee that the prediction does not change

```
Verify whether for all admissible perturbed graphs \tilde{G}:

argmax_{class_y} f(\tilde{G})_y \stackrel{?}{=} argmax_{f(G)_y} f(G)_y
perturbed graph - clean graph
```











(LOWER BOUND ON) THE WORST-CASE MARGIN



COMPUTING THE WORST-CASE MARGIN

Model-specific Certificates

Model-agnostic Certificates

PAGERANK-BASED MODELS

Predictions are a linear function of personalized PageRank π_G

$$f(G)_{y} = \pi_{G}^{T} H_{:,y}$$
personalized PageRank
vector for a given node
- "logits" for class y

Models in this family: **PPNP**, PushNet, Label/Feature Propagation

BACKGROUND: PERSONALIZED PAGERANK (PPR)

Stationary distribution of a random walk with teleport



The teleport probability α controls the effective neighborhood size

BACKGROUND: PPR AND LABEL PROPAGATION

Repeatedly diffuse inital "beliefs" using the graph

$$H^{(0)} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \\ \vdots & \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 0 \\ H^{(t+1)} = (1 - \alpha)D^{-1}AH^{(t)} + \alpha H^{(0)} \\ H^{(\infty)} = \Pi^{\text{ppr}}H^{(0)}$$

Probability that some node belongs to class $y \propto \pi_G^T H_{:,y}$

(A)PPNP: PREDICT THEN PROPAGATE

First map node features to initial beliefs then diffuse them with PPR



PAGERANK OPTIMIZATION



$$m^* = \min_{\substack{\tilde{G} \in \text{admisible} \\ \text{perturbed graphs}}} \min_{\substack{\text{class } y \neq y^* \\ y \neq y^*}} \pi_{\tilde{G}}^T \left(H_{:,y^*} - H_{:,y} \right)$$

Minimize a linear function of PageRank over the set of graphs

THREAT MODEL – WHAT IS ADMISSIBLE?

An attacker has control over a set \mathcal{F} of fragile edges



General and flexible:

Scenario I: $\mathcal{F} = \mathcal{E}$ # remove edgesScenario 2: $\mathcal{F} = (\mathcal{V} \times \mathcal{V}) \setminus \mathcal{E}$ # add edgesScenario 3: $\mathcal{F} = \dots$

Global budget: perturb up to B edges in total Local budget: perturb up to b_v edges for node v

PAGERANK OPTIMIZATION





PAGERANK OPTIMIZATION





EXACT CERTIFICATES FOR PAGERANK-BASED MODELS





- Feature Propagation
- Label Propagation

EXACT CERTIFICATES FOR PAGERANK-BASED MODELS



ROBUST TRAINING



COMPUTING THE WORST-CASE MARGIN

Model-specific Certificates

Model-agnostic Certificates

+ Takes advantage of problem structure (exact guarantees)

Limited applicability

RANDOMLY SMOOTHED CLASSIFIERS

Any base classifier $f: \mathcal{X} \to \mathcal{Y}$ Randomization scheme $\phi(\mathbf{x})$

Certify a smoothed classifier g

$$g(\mathbf{x}) = \underset{y \in \mathcal{Y}}{\operatorname{argmax}} \Pr(f(\boldsymbol{\phi}(\mathbf{x})) = y)$$

majority vote y*



RANDOMLY SMOOTHED CLASSIFIERS

Guarantee that the majority does not change in a ball \mathcal{B}_r around \boldsymbol{x}



RANDOMLY SMOOTHED CLASSIFIERS

Guarantee that the majority does not change in a ball \mathcal{B}_r around \boldsymbol{x}

$$f(\mathbf{x}) = \mathbf{0}$$
, but $g(\mathbf{x}) = \mathbf{0}$

Verify whether for all $\widetilde{x} \in \mathcal{B}_r(x)$ $\Pr(f(\phi(\widetilde{x})) = \bullet) \stackrel{?}{>} 0.5$



TWO NECESSARY IMPROVEMENTS

I. Sparsity-aware smoothing



TWO NECESSARY IMPROVEMENTS

I. Sparsity-aware smoothing



2. Dramatically more efficient: $O(n^4) \rightarrow O(r) \triangleleft ---$ certified radius r



DERIVING THE CERTIFICATE

The smoothed classifier is certifiably robust if

min
$$\Pr(f(\boldsymbol{\phi}(\widetilde{\boldsymbol{x}})) = y^*) \stackrel{?}{>} 0.5$$

subject to:

 $\widetilde{x} \in \mathcal{B}_r(x) \leftarrow -$ admissible

Find the \widetilde{x} that minimizes the probability of the majority vote y^*

CONSTANT LIKELIHOOD RATIO REGIONS



GNNS HAVE DIFFERENT ROBUSTNESS TRADE-OFFS





GNNS HAVE DIFFERENT ROBUSTNESS TRADE-OFFS

GRAPH-LEVEL CLASSIFICATION

CERTIFYING IMAGENET

| Certificate | Туре | Time | <i>r</i> = 1 | r = 3 | <i>r</i> = 5 | r = 7 |
|-------------------------|------------|----------|--------------|-------|--------------|-------|
| Cohen et al. (2019) | Continuous | < I sec. | 0.372 | 0.226 | 0.170 | 0.138 |
| Dvijotham et al. (2020) | Discrete | < I sec. | 0.362 | 0.224 | 0.136 | 0 |
| Lee et al. (2019) | Discrete | 4 days | 0.538 | 0.338 | 0.244 | 0.176 |
| Ours | Discrete | < I sec. | 0.538 | 0.338 | 0.244 | 0.176 |

COMPUTING THE WORST-CASE MARGIN

Model-specific Certificate

+ Takes advantage of problem structure (exact guarantees)

Model-agnostic Certificate

+ Can be applied to any classifier for discrete data (all GNNs)

- Limited applicability

Does not capture all properties of the classifier

WHY ARE THE CERTIFICATES SO PESSIMISTIC?

They treat each node independently

COLLECTIVE ROBUSTNESS CERTIFICATES

Verify how many nodes are simultaneously robust

FUSE MULTIPLE SINGLE-NODE CERTIFICATES

Into a **provably stronger** certificate by exploiting interdependence

Budget allocation problem: pick k edges to misclassify most nodes

COLLECTIVE = DRAMATIC IMPROVEMENTS

MORE REALISTIC THREAT MODELS

SUMMARY

GNNs are vulnerable to small perturbations \Rightarrow Verify robustness

Model-specific: Exact guarantees for a family of PageRank-based models Model-agnostic: Turn any classifier into a certifiable smoothed classifier Collective: Certify simultaneously robust nodes (dramatic improvement)

Open: knowledge graphs, code, graph-level, representations, poisoning, ...

www.daml.in.tum.de/graph-cert www.daml.in.tum.de/sparse-smoothing www.daml.in.tum.de/collective-robustness

