### Reliable Graph Machine Learning

Simon Geisler



Reliability

## This talk focuses on reliability in the sense of worst case robustness w.r.t. perturbations of the input – commonly referred to as **adversarial robustness**.



Adversarial Robustness

"pig"

#### **Predicted class:**



7



م Adversarial Robustness

#### **Predicted class:**



#### "airliner"





Adversarial Robustness



ML systems classify the adversarially modified STOP sign as a speed limit sign, [Eykholt 2018]

Kevin Eykholt, et al. "Robust physical-world attacks on deep learning visual classification." CVPR 2018.



### Motivations of Adversarial Robustness

- security:We do not want a real-world adversary to succeed- safety:Our application entails requirements on worst-case<br/>behaviour of  $f_{\theta}$ 

local sensitivity: We have some expectation about local behavior of f<sub>θ</sub>
 generalization: The admissible perturbations B are linked to the semantics of the input instance









Adversarial attacks are *typically* approximate: **Upper bound** 

Certificates *typically* require some relaxation: **Lower bound** 

Determining the **true robustness** is *typically* intractable



# $\underbrace{}_{\mathbf{A}} \text{Perspectives on Adversarial Robustness}}_{\mathbf{B}} \in \mathcal{B}(\mathbf{A}) \text{ s.t. } \operatorname{argmax}(f_{\theta}(\mathbf{A})) \neq \operatorname{argmax}(f_{\theta}(\mathbf{A}))$

#### **Assessing Robustness**

- Upper bound: Attacks
- Lower bound: Certification (will not be covered today)

#### Improving Robustness

- Invariant model architectures
- Enhanced model architecture
- Robust/adversarial training





#### I. Adversarial Robustness

- 2. Adversarial Attacks (Assessing Robustness)
  - a. Application: Combinatorial Optimization
  - b. Scalable Attacks for Structure Perturbations
  - c. Loss Functions for Interconnected Predictions: Node Classification
- 3. Improving Robustness
  - a. Invariant model architectures
  - b. Enhanced model architectures
  - c. Robust/adversarial training





- I. Adversarial Robustness
- 2. Adversarial Attacks (Assessing Robustness)

#### a. Application: Combinatorial Optimization

- b. Scalable Attacks for Structure Perturbations
- c. Loss Functions for Interconnected Predictions: Node Classification
- 3. Improving Robustness
  - a. Invariant model architectures
  - b. Enhanced model architectures
  - c. Robust/adversarial training



### <u>Combinatorial</u> Optimization: TSP

From a problem instance, find a perturbed problem instance  $\tilde{x}$  that maximizes the surrogate loss  $L(f_{\theta}(\tilde{x}), \tilde{Y})$ .



Original problem instance  $\pmb{x}$  Perturbed problem instance  $\widetilde{\pmb{x}}$ 

#### $\rightarrow$ The admissible perturbations are linked to the semantics

Geisler et al. "Generalization of Neural Combinatorial Solvers Through the Lens of Adversarial Robustness." ICLR 2022.

Simon Geisler



SAT Perturbation Model  

$$\overbrace{(l_1 \vee \neg l_2)}^{c_1} \wedge \overbrace{(l_2)}^{c_2} \wedge \overbrace{(\neg l_1 \vee l_2 \vee l_3)}^{c_3} \wedge \overbrace{(l_1 \vee l_3)}^{c_4} \wedge \overbrace{(\neg l_2 \vee \neg l_3)}^{c_5}$$

$$= \overbrace{(1 \vee 0)}^{c_1} \wedge \overbrace{(0 \vee 1 \vee 0)}^{c_2} \wedge \overbrace{(1 \vee 0)}^{c_3} \wedge \overbrace{(0 \vee 1 \neg )}^{c_4} = 1$$

*Invariance:* we can insert/remove literals as long as a literal of truth assignment *Y* remains in each clause

$$\underbrace{\overbrace{(l_1 \lor \neg l_2)}^{c_1} \land \underbrace{\overbrace{\neg l_1 \lor l_2}^{c_2}}_{\text{insert}} \land \underbrace{\overbrace{(\neg l_1 \lor l_2}^{c_3} \lor \underbrace{\lor l_3}_{\text{remove}}) \land \underbrace{\overbrace{(l_1 \lor l_3)}^{c_4} \land \underbrace{(\neg l_2 \lor \neg l_3)}_{\text{remove}}$$



Geisler et al. "Generalization of Neural Combinatorial Solvers Through the Lens of Adversarial Robustness." ICLR 2022.





 $rac{1}{7}$  Changing approx. 0.5% of the literals suffices to push the accuracy below 50%  $rac{1}{7}$ 

Geisler et al. "Generalization of Neural Combinatorial Solvers Through the Lens of Adversarial Robustness." ICLR 2022.

Simon Geisler





I. Adversarial Robustness

### 2. Adversarial Attacks (Assessing Robustness)

a. Application: Combinatorial Optimization

#### b. Scalable Attacks for Structure Perturbations

- c. Loss Functions for Interconnected Predictions: Node Classification
- 3. Improving Robustness
  - a. Enhanced model architectures
  - b. Invariant model architectures
  - c. Robust/adversarial training



### Adversarial Attack for Structure Perturbations

Adversarial attack on a fixed GNN  $f_{\theta}$  with loss  $\mathcal{L}$  and budget  $\Delta$ :

- $\max \mathcal{L}(f_{\theta}( \bullet)) \text{ s.t. } \| \bullet \bullet \|_{0} \leq \Delta$
- $\stackrel{\checkmark}{\sim}$  Results in a discrete and non-convex optimization problem
- A graph with n nodes has  $n^2$  possible edges

 $\rightarrow$  Projected Randomized Block Coordinate Descent (PR-BCD) (complexity  $O(\Delta)$ )





Adversarial Attacks on Graph Neural Networks

### $\max_{\mathbf{P}} \mathcal{L}(f_{\theta}(\mathbf{A} \bigoplus \mathbf{P}, \mathbf{X})) \text{ s.t. } \sum_{\mathbf{P}} \leq \Delta \land \mathbf{P} \in \{0, 1\}^{n \times n} \land \mathbf{P} = \mathbf{P}^{\mathrm{T}}$

For the applicability of gradient-based optimization

- (1)  $\mathbf{P} \in \{0,1\}^{n \times n}$  is relaxed to  $\mathbf{P} \in [0,1]^{n \times n}$  during the attack
- (2) Each entry in  $\mathbf{P} \in [0,1]^{n \times n}$  represents the **probability of an edge flip**
- (3) After the attack, we sample from  $\mathbf{P} \in [0,1]^{n \times n}$  to obtain  $\mathbf{P} \in \{0,1\}^{n \times n}$

**Assumption:** model can handle edge weights





#### $\rightarrow$ Effectively relaxes the constraint from $\|\cdot\|_0$ to $\|\cdot\|_1$



### Adversarial Attacks on Graph Neural Networks

### $\max_{\mathbf{P}} \mathcal{L}(f_{\theta}(\mathbf{A} \bigoplus \mathbf{P}, \mathbf{X})) \text{ s.t. } \sum_{\mathbf{P}} \leq \Delta \land \mathbf{P} \in \{0, 1\}^{n \times n} \land \mathbf{P} = \mathbf{P}^{\mathrm{T}}$



 $\checkmark$  Number of "parameters" scales with  $\mathcal{O}(n^2)$  with number of nodes n.





Our variant of Projected Randomized Block Coordinate Descent (PR-BCD) that maintains a sparse parameter space throughout the optimization



 $\max_{\mathbf{P}} \mathcal{L}(f_{\theta}(\mathbf{A} \bigoplus \mathbf{P}, \mathbf{X})) \text{ s.t. } \sum_{\mathbf{P}} \mathbf{P} \leq \Delta \land \mathbf{P} \in \{0, 1\}^{n \times n}$ 

Example: Undirected Graph with n=5 nodes and budget  $\Delta=2$ 

Our Proposed Attack

#### Our variant of Projected Randomized Block Coordinate Descent (PR-BCD) that maintains a *sparse parameter space* throughout the optimization



#### $\max_{\mathbf{P}} \mathcal{L}(f_{\theta}(\mathbf{A} \bigoplus \mathbf{P}, \mathbf{X})) \text{ s.t. } \sum_{\mathbf{P}} \mathbf{P} \leq \Delta \land \mathbf{P} \in \{0, 1\}^{n \times n}$

Example: Undirected Graph with n=5 nodes and budget  $\Delta=2$ 

Simon Geisler



### <u>Scalable Attacks for Structure Perturbations</u>

Dataset	# Nodes <b>n</b>	Size (dense)	Size (sparse)	
Citeseer	2.1k	17.8MB	146.7kB	
Cora ML	2.8k	31.9MB	319.2kB	<ul> <li>Previous work</li> </ul>
PubMed	19.7k	I.6GB	I.8MB	
arXiv *	169.3k	114.7GB	23.3MB	
Products	2.5M	24.0TB	2.5GB	- Ours (first-order attack
Papers	III.IM	49.3PB	32.3 I GB	``

#### $\rightarrow$ We are the first to study adversarial robustness on such massive graphs.

\* There are a few exceptions that study adversarial robustness for some special cases on graphs with similar size



### ◇QI: Are GNNs robust if applied to large graphs?

Adversarial attack on GCN on Products (2.5 mio. nodes, 124 mio. edges):





### PR-BCD is available in PyTorch Geometric

from torch\_geometric.contrib.nn import PRBCDAttack





What are reasonable / semantic preserving / ... perturbations  $\leftarrow \in \mathcal{B}(\leftarrow)$ ?

 $\rightarrow$  Open question and most certainly application specific

It is reasonable to design flexible attacks, e.g., also supporting local constraints:

**Global constraint** 
$$\sum_{i,j} \mathbf{P}_{i,j} \leq \Delta^{(g)}$$
 and **local constraints**<sup>\*</sup>  $\sum_{j} \mathbf{P}_{i,j} \leq \Delta_{i}^{(l)}$ 

In Geisler et al. "Adversarial Training for Graph Neural Networks" arXiv 2023 we extend the presented attack to also obey local constraints:

Geisler et al. "Adversarial Training for Graph Neural Networks." In arXiv. 2023.

 $\ast$  with symmetrized **P** 

Simon Geisler







I. Adversarial Robustness

### 2. Adversarial Attacks (Assessing Robustness)

- a. Application: Combinatorial Optimization
- b. Scalable Attacks for Structure Perturbations

### c. Loss Functions for Interconnected Predictions: Node Classification

- 3. Improving Robustness
  - a. Enhanced model architectures
  - b. Invariant model architectures
  - c. Robust/adversarial training



PR-BCD is an efficient first-order optimization method to approximate

$$\max \mathcal{L}(f_{\theta}( \bullet)) \text{ s.t. } \| \bullet \bullet \bullet \bullet \|_{0} \leq \Delta$$
  
Surrogate losses for interconnected targets

... but do we optimize for a sensible target with common losses?



### Surrogate Losses for "Global" Attacks



Simon Geisler



### Surrogate Losses for "Global" Attacks (binary)



Logit margin  $\max_{c \neq c^*} z_{c^*} - z_c$ 





- I. Adversarial Robustness
- 2. Adversarial Attacks (Assessing Robustness)
  - a. Application: Combinatorial Optimization
  - b. Scalable Attacks for Structure Perturbations
  - c. Loss Functions for Interconnected Predictions: Node Classification

### 3. Improving Robustness

- a. Invariant model architectures
- b. Enhanced model architectures
- c. Robust/adversarial training



### How Do Transformers Encode Code?

```
def f1_score(pred, label):
    correct = pred == label
    tp = (correct & label).sum()
    fn = (~correct & pred).sum()
    fp = (~correct & ~pred).sum()
    precision = tp / (tp + fp)
    recall = tp / (tp + fn)
    return (
        2 * (recall * precision) /
        (recall * precision)
    )
```

#### Tokenization:





### Motivation for Directed Graphs: Encode symmetries



 $\rightarrow$  The model itself and the modelling of the input data greatly impact robustness

Geisler et al. "Transformers Meet Directed Graphs." ICML 2023.

Reliable Graph Machine Learning



**Note:** Our construction via the

Abstract Syntax Tree (AST)

covers loops, exceptions, etc.



- I. Adversarial Robustness
- 2. Adversarial Attacks (Assessing Robustness)
  - a. Application: Combinatorial Optimization
  - b. Scalable Attacks for Structure Perturbations
  - c. Loss Functions for Interconnected Predictions: Node Classification

#### 3. Improving Robustness

- a. Invariant model architectures
- b. Enhanced model architectures
- c. Robust/adversarial training



### Scalable Robust Message Passing Aggregation

$$\mathsf{E.g., a GCN: } \mathbf{h}_{v}^{(l)} \approx \sigma^{(l)} \left( \mathsf{MEAN} \left\{ \left( \mathsf{A}_{vu}, \ \mathbf{h}_{u}^{(l-1)} \, \mathsf{W}^{(l)} \right), \forall \, u \in \mathcal{N}(v) \cup v \right\} \right)$$

Geisler et al.''Robustness of Graph Neural Networks at Scale.'' NeurIPS 2021. Geisler et al.''Reliable Graph Neural Networks via Robust Aggregation.'' NeurIPS 2020.


Geisler et al. "Robustness of Graph Neural Networks at Scale." NeurIPS 2021. Geisler et al. "Reliable Graph Neural Networks via Robust Aggregation." NeurIPS 2020.

Scalable Robust Message Passing Aggregation  
Our Soft Median: 
$$\mu(\mathbf{X}) = \operatorname{softmax} \left(\frac{-\mathbf{d}}{T}\right)^T \mathbf{X} = \mathbf{s}^T \mathbf{X}$$
 ≈  $\operatorname{argmin}_{\mathbf{x}' \in \mathbb{X}} \|\bar{\mathbf{x}} - \mathbf{x}'\|$ 

Distance to dim.-wise Median:  $\mathbf{d}_i = \|\mathbf{\bar{x}} - \mathbf{x}_i\|$  Input Embeddings:  $\mathbf{X}, \mathbf{X}$ 



Geisler et al. "Robustness of Graph Neural Networks at Scale." NeurIPS 2021. Geisler et al. "Reliable Graph Neural Networks via Robust Aggregation." NeurIPS 2020.

$$\frac{1}{\sqrt{2}} \frac{\text{Scalable Robust Message Passing Aggregation}}{\text{Our Soft Median:}} \mu(\mathbf{X}) = \operatorname{softmax} \left(\frac{-\mathbf{d}}{T}\right)^{\mathsf{T}} \mathbf{X} = \mathbf{s}^{\mathsf{T}} \mathbf{X} \quad \approx \operatorname{argmin}_{\mathbf{x}' \in \mathbb{X}} \|\bar{\mathbf{x}} - \mathbf{x}'\|$$

Distance to dim.-wise Median:  $\mathbf{d}_i = \|\overline{\mathbf{x}} - \mathbf{x}_i\|$  Input Embeddings: **X**, **X** 





# Evaluation Pitfall: Non-Adaptive Attacks

What we should evaluate for model  $f_{\theta}$ :



What is being evaluated instead (with surrogate model  $g_{\theta} \neq f_{\theta}$ ):





52

How to Design Adaptive Attacks?

Many defenses  $f_{\theta}$  disallow certain connections (e.g. Jaccard GCN, SVD GCN)



Geisler et al. "Are Defenses for Graph Neural Networks Robust?" NeurIPS 2022.



## Adaptive vs. Non-Adaptive Attacks



Geisler et al. "Are Defenses for Graph Neural Networks Robust?" NeurIPS 2022.



## Adaptive vs. Non-Adaptive Attacks



Geisler et al. "Are Defenses for Graph Neural Networks Robust?" NeurIPS 2022.



## Adaptive vs. Non-Adaptive Attacks



#### $\rightarrow$ Our defense motivated by robust statistics seems to be the only effective option





- I. Adversarial Robustness
- 2. Adversarial Attacks (Assessing Robustness)
  - a. Application: Combinatorial Optimization
  - b. Scalable Attacks for Structure Perturbations
  - c. Loss Functions for Interconnected Predictions: Node Classification

#### 3. Improving Robustness

- a. Invariant model architectures
- b. Enhanced model architectures

### c. Robust/adversarial training



# Adversarial Training in the Inductive Setting

We adapt the common benchmark setup for node-level predictions:

- 5% of nodes are used for training
- 80% of nodes are unlabeled
- 5% of nodes are used for validation
- 10% of nodes are used for test

NOT included in training graph

→ Semi-supervised **inductive learning** 

Geisler et al. "Adversarial Training for Graph Neural Networks." In arXiv. 2023.



# Adversarial Training in the Inductive Setting

<u>Regular training objective</u> where we sample the graph extended by test nodes and test labels  $(\mathcal{G}', y') \sim \mathcal{D}(\mathcal{G}, y)$  with test set  $\mathbb{I}_{\text{Test}}$ :

$$\min_{\theta} \mathbb{E}_{(\mathcal{G}', y') \sim \mathcal{D}(\mathcal{G}, y)} \sum_{i \in \mathbb{I}_{\text{Test}}} \mathcal{L}(f_{\theta}(\mathcal{G}')_{i}, y_{i}')$$

Adversarial training objective:

$$\min_{\theta} \mathbb{E}_{\substack{(\mathcal{G}', y') \sim \mathcal{D}(\mathcal{G}, y)}} \max_{\tilde{\mathcal{G}}' \in \mathcal{B}(\mathcal{G}')} \sum_{i \in \mathbb{I}_{\text{Test}}} \mathcal{L}\left(f_{\theta}\left(\tilde{\mathcal{G}}'\right)_{i}, y_{i}'\right)$$

Geisler et al. "Adversarial Training for Graph Neural Networks." In arXiv. 2023.



Geisler et al. "Adversarial Training for Graph Neural Networks." In arXiv. 2023.



# Model Choice for Adversarial Training

We propose to use flexible GNN models that fit in the framework

$$\sum_{k=0}^{K} \gamma_k \boldsymbol{L}^k \boldsymbol{H} = \boldsymbol{V} \operatorname{diag}(\hat{g}(\boldsymbol{\Lambda})) \boldsymbol{\overline{V}}^T \boldsymbol{H} = \boldsymbol{V} \begin{bmatrix} \hat{g}(\lambda_1) & 0 & \cdots & 0 \\ 0 & \hat{g}(\lambda_2) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \hat{g}(\lambda_n) \end{bmatrix} \boldsymbol{\overline{V}}^T \boldsymbol{H}$$

A sufficiently flexible message passing scheme (we use K = 10)
 We can interpret the learned message passing  $\sum_{k=0}^{K} \gamma_k L^k$ 

Geisler et al. "Adversarial Training for Graph Neural Networks." In arXiv. 2023.



Robust Message Passing CharacteristicsModel: $\sum_{k=0}^{K} \gamma_k L^k H = V \operatorname{diag}(\hat{g}(\Lambda)) \overline{V}^T H$ Setup: Citeseer with 20%<br/>of edges perturbed

#### Polynomial interpretation:

Spectral filter  $\widehat{g}(oldsymbol{\Lambda})$  :





### Efficacy of Adversarial Training

Adversarial attacks: LR-BCD: w/ local constraints PR-BCD: w/o local constraints

Dataset: Citeseer

Model	Adv.	A. eval. $\rightarrow$	Classe	Clean $LR$ -BCD PR-BCD $\epsilon=0.1$		LR-BCD	PR-BCD	Certifiable accuracy / sparse smoothing		
	trn.	A. trn. $\downarrow$	Clean			€=0.25		Clean	3 add.	5 del.
GCN	×	-	$72.0\pm2.5$	$54.7\pm2.8$	$51.7\pm2.8$	$45.3\pm3.4$	$38.0\pm3.8$	$38.3 \pm 11.5$	$1.7\pm0.7$	$4.8\pm1.5$
GAT	X	-	$-3.6\pm2.7$	$-3.9\pm3.4$	$+0.5\pm3.5$	$-15.9\pm5.3$	$-2.3\pm7.3$	$-14.0\pm12.0$	$-1.7\pm0.7$	$-4.8 \pm 1.5$
APPNP	X	-	$+0.2 \pm 1.1$	$+1.7\pm0.7$	$+1.9 \pm 1.4$	$+3.0 \pm 1.2$	$+2.2 \pm 2.5$	$+8.9\pm9.1$	$+31.2 \pm 6.4$	$+31.6 \pm 6.4$
GPRGNN	X	-	$+2.2 \pm 4.3$	$+4.2\pm2.7$	$+3.6\pm4.9$	$+5.5 \pm 3.9$	$+7.9\pm4.6$	$+17.9\pm6.9$	+42.4 $\pm$ 4.4	$+41.3 \pm 3.7$
ChebNetII	X	-	$+1.1 \pm 2.2$	$+5.8\pm2.5$	$+5.0\pm2.4$	$+10.4\pm2.6$	$+7.6 \pm 3.4$	$+24.6\pm9.9$	$+55.6 \pm 1.2$	$+54.0\pm0.9$
SoftMedian	X	-	+0.9 $\pm$ 1.7	$+9.5\pm2.2$	$+9.3\pm1.9$	$+16.2 \pm 2.4$	$+14.6\pm2.9$	$+25.2\pm10.5$	$+60.3 \pm 1.4$	$+57.5\pm0.8$
GCN	$\checkmark$	LR-BCD	$-0.2 \pm 1.2$	$+7.8\pm1.6$	$+5.9\pm1.5$	$+10.9\pm2.1$	$+8.1\pm2.3$	$-3.0\pm9.4$	$+10.7\pm4.6$	$+13.1 \pm 5.1$
		PR-BCD	+ $0.0 \pm 1.9$	+6.9 $\pm$ 0.9	$+5.3 \pm 1.6$	$+8.6 \pm 2.3$	$+5.8\pm2.4$	$+4.2 \pm 14.6$	$+10.1 \pm 5.5$	$+11.4 \pm 4.5$
GAT	$\checkmark$	LR-BCD	+ $0.8\pm1.6$	$+5.9\pm3.7$	$+9.0\pm3.0$	$+8.4 \pm 5.1$	$+13.1 \pm 3.5$	$-2.0\pm23.0$	$+1.4 \pm 1.7$	+4.0 $\pm$ 2.1
		PR-BCD	$+1.1 \pm 2.2$	$+8.9\pm2.8$	$+13.2 \pm 3.7$	$+10.3 \pm 2.3$	$+21.8\pm4.8$	$-10.1 \pm 13.7$	$+1.2 \pm 2.0$	$+2.8 \pm 1.2$
APPNP	/	LR-BCD	$-0.8\pm1.3$	$+7.6 \pm 1.6$	$+5.6 \pm 1.9$	$+11.1 \pm 1.9$	$+8.3\pm3.6$	$+21.7\pm9.5$	+41.3 $\pm$ 2.9	$+44.1 \pm 1.2$
	V	PR-BCD	$-0.2 \pm 2.2$	$+6.2 \pm 2.3$	$+5.5 \pm 3.1$	$+9.0\pm2.9$	$+6.5 \pm 3.8$	$+19.3 \pm 7.2$	$+41.0 \pm 3.8$	$+41.9 \pm 3.0$
GPRGNN	$\checkmark$	LR-BCD	$+1.7 \pm 3.0$	+15.7 $\pm$ 3.6	$+13.6 \pm 3.5$	+24.8 $\pm$ 4.2	$+22.9\pm4.3$	+34.0 $\pm$ 11.5	$+67.4 \pm 1.5$	$+65.0\pm2.5$
		PR-BCD	$+0.6\pm3.6$	$+15.0 \pm 3.6$	+15.9 $\pm$ 4.2	$+23.2 \pm 4.3$	+26.3 $\pm$ 5.4	$+32.9 \pm 11.8$	+69.0 $\pm$ 1.5	+65.9 $\pm$ 2.3
ChebNetII	$\checkmark$	LR-BCD	$+3.7\pm1.4$	$+11.5 \pm 1.6$	$+11.5 \pm 1.1$	$+16.4 \pm 1.9$	$+16.0\pm2.9$	$+31.5 \pm 12.7$	+62.3 $\pm$ 1.9	$+59.8\pm2.5$
		PR-BCD	$+3.4 \pm 1.7$	$+13.2\pm2.2$	$+14.3\pm1.2$	$+19.5\pm1.6$	$+19.8\pm2.3$	$+30.7\pm13.2$	$+65.4\pm2.6$	$+62.9\pm3.7$

Geisler et al. "Adversarial Training for Graph Neural Networks." In arXiv. 2023.





- I. Adversarial Robustness
- 2. Adversarial Attacks (Assessing Robustness)
  - a. Application: Combinatorial Optimization
  - b. Scalable Attacks for Structure Perturbations
  - c. Loss Functions for Interconnected Predictions: Node Classification
- 3. Improving Robustness
  - a. Invariant model architectures
  - b. Enhanced model architectures
  - c. Robust/adversarial training



Summary

Email:geisler@in.tum.deTwitter/X:@geisler\_si

#### **Covered works / references:**

Geisler et al. "Adversarial Training for Graph Neural Networks." arXiv 2023.

Geisler et al. "Transformers Meet Directed Graphs." ICML 2023.

Geisler et al. "Are Defenses for Graph Neural Networks Robust?" NeurIPS 2022.

Geisler et al. "Generalization of Neural Combinatorial Solvers Through the Lens of Adversarial Robustness." ICLR 2022.

Geisler et al. "Robustness of Graph Neural Networks at Scale." NeurIPS 2021.

Geisler et al. "Reliable Graph Neural Networks via Robust Aggregation." NeurIPS 2020.





- I. Adversarial Robustness
- 2. Graph Neural Networks
- 3. Assessing Robustness
  - a. Certification
  - b. Adversarial Attacks
- 4. Improving Robustness
  - a. Enhanced model architecture
  - b. Invariant model architectures
  - c. Robust/adversarial training





### Classes of Certificates

#### Strategy I:

"White-box" certificates

- + Take advantage of the problem structure for accurate certificates,
  i.e. tight lower bounds
- Special derivation for every model required, i.e. expensive development

### Strategy II: "Black-box" certificates

- + Can be applied to any classifier for discrete data (including all GNNs) out of the box
- Bounds might be less tight since model structure not exploited

#### Is there something in between? Strategy III: "Gray-box" certificates

Zügner et al. Certifiable Robustness of Graph Convolutional Networks under Structure Perturbations. KDD 2020 Bojchevski at al. Certifiable Robustness to Graph Perturbations. NeurIPS 2019 Zügner et al. Certifiable Robustness and Robust Training for Graph Convolutional Networks. KDD 2019 Schuchardt et al. Localized Randomized Smoothing for Collective Robustness Certification. ICLR 2023 Schuchardt et al. Collective Robustness Certificates: Exploiting Interdependence in GNNs. ICLR 2021 Bojchevski et al. Efficient Robustness Certificates for Discrete Data. ICML 2020





- Enhance black-box by general model properties
  - message-passing princple
  - invariances
  - . . .



Adversarial message



Schuchardt, Günnemann. Invariance-Aware Randomized Smoothing Certificates. NeurIPS 2022 Scholten, Schuchardt, Geisler, Bojchevski, Günnemann. Randomized Message-Interception Smoothing: Gray-box Certificates for Graph Neural Networks. NeurIPS 2022



- Enhance black-box by general model properties
  - message-passing princple
  - invariances



Schuchardt, Günnemann. Invariance-Aware Randomized Smoothing Certificates. NeurIPS 2022

Scholten, Schuchardt, Geisler, Bojchevski, Günnemann. Randomized Message-Interception Smoothing: Gray-box Certificates for Graph Neural Networks. NeurIPS 2022

### Individual vs. Collective Certificates

- Often: classifier that simultaneously outputs multiple predictions give a single input
  - node classification, image segmentation, named-entity recognition
- Problem: "Standard" certificates care about a single prediction only
- Collective robustness certificate → number of predictions that are simultaneously guaranteed to remain stable







#### Task: Node-Level Classification







#### Transformers Meet Directed Graphs Previous Transformer **Graph Transformer** Ours How to generalize transformers to directed graphs? $\rightarrow$ Direction-aware positional encodings: Jndirected Sequences Graphs Spectral encodings: Magnetic Laplacian a) b) Directional random walks Space of directed Graphs

- 2. Directed graphs are important for many applications
- $\rightarrow$  Appropriate modeling of input modalities:
  - If ignoring the direction, the model might sacrifice expressivity a)

We can leverage symmetries in the input domain















d- # of features/<br/>embedding dim. $2\pi 10,000$ - largest wave lengthv- node index/positionj- frequency indexx- signal (spatial domain) $\hat{x}$ - signal (spectral domain)n- sequence length

• Sinusoidal Positional Encodings  $\mathrm{PE}^{(\sin)} \in \mathbb{R}^{n \times d}$ 

$$PE_{v,j+d/2}^{(sin)} = sin(v \cdot 1/10, 000^{2j/d})$$
$$PE_{v,j}^{(sin)} = cos(v \cdot 1/10, 000^{2j/d})$$

• Positional encodings inspired by Discrete Fourier Transformation (DFT):

$$\hat{x}_j = \sum_{\nu=0}^{n-1} x_{\nu} \left[ \cos\left(\nu \cdot \frac{2\pi}{n}j\right) - i \cdot \sin\left(\nu \cdot \frac{2\pi}{n}j\right) \right] = \sum_{\nu=0}^{n-1} x_{\nu} e^{\frac{-i 2\pi\nu j}{n}}$$

### $\sim$ Sinusoidal Encodings $\leftrightarrow$ Laplacian Eigenvectors



#### Assumption: graphs are undirected

Graph Fourier Transformation: Eigenvectors of combinatorial/graph Laplacian







Discrete

Reliable Graph Machine Learning

related to

Sinusoidal

Circular Convolution:

= IDFT(DFT( $h(\omega) * x$ )) y  $= \text{IDFT}\left(\hat{h}(\omega) \cdot \text{DFT}(\boldsymbol{x})\right)$  $= 1/n \mathbf{S} \operatorname{diag}(\hat{h}(\omega)) \overline{\mathbf{S}}^T \mathbf{x}$ with DFT  $\hat{x} = \overline{S}^T x$ ,  $\mathsf{IDFT} \mathbf{y} = S\widehat{\mathbf{y}},$ Fourier sinusoids  $S \in \mathbb{C}^{n \times n}$ , and diagonal filter diag $(\hat{h}(\omega)) \in \mathbb{R}^{n \times n}$ .

$$\overline{\mathbf{S}}^{T} = \begin{bmatrix} 1 & 1 & 1 & \cdots & 1 \\ 1 & \omega & \omega^{2} & \cdots & \omega^{n-1} \\ 1 & \omega^{2} & \omega^{4} & \cdots & \omega^{2(n-1)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & \omega^{n-1} & \omega^{2(n-1)} & \cdots & \omega^{(n-1)^{2}} \end{bmatrix}$$

$$\operatorname{diag}(\hat{h}(\omega)) = \begin{bmatrix} h(\omega^0) & 0 & \cdots & 0 \\ 0 & h(\omega^1) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & h(\omega^{n-1}) \end{bmatrix}$$

with 
$$\omega = e^{\frac{-i 2\pi v}{n}}$$

Signal Processing: DFT ↔ GFT

Circular Convolution:

= IDFT(DFT( $h(\omega) * \mathbf{x}$ )) y  $= \text{IDFT}\left(\hat{h}(\omega) \cdot \text{DFT}(\boldsymbol{x})\right)$  $= 1/n \mathbf{S} \operatorname{diag}(\hat{h}(\omega)) \overline{\mathbf{S}}^T \mathbf{x}$ with DFT  $\hat{x} = \overline{S}^T x$ .  $\mathsf{IDFT} \mathbf{y} = S\widehat{\mathbf{y}},$ Fourier sinusoids  $S \in \mathbb{C}^{n \times n}$ , and diagonal filter diag $(\hat{h}(\omega)) \in \mathbb{R}^{n \times n}$ . Convolution on graphs:

$$= \operatorname{IGFT}(\operatorname{GFT}(g(\Lambda) * x))$$
  
= IGFT( $\widehat{g}(\Lambda) \cdot \operatorname{GFT}(x)$ )  
=  $V \operatorname{diag}(\widehat{g}(\Lambda)) \overline{V}^T x$ 

with eigendecomposition  $L = V\Lambda \overline{V}^T$ of graph Laplacian L = D - Awith adjacency matrix  $A \in \mathbb{R}^{n \times n}_{\geq 0}$ and degree matrix  $D = diag(A\vec{1})$ 



#### Symmetrization — real and symmetric Laplacian



### Magnetic Laplacian

$$\boldsymbol{L}^{(q)} := \boldsymbol{D}_{s} - \boldsymbol{A}_{s} \odot \exp(i\boldsymbol{\Theta}^{(q)})$$

$$\Theta_{u,v}^{(q)} := 2\pi q \left( A_{u,v} - A_{v,u} \right)$$

 $A_{s} \in \mathbb{R}_{\geq 0}^{n \times n}$  $D_{s}$ q $exp/\bigcirc$  $i = \sqrt{-1}$ 

symmetrized adjacency matrix diagonalized degree matrix potential *(hyperparameter)* element-wise exp./product complex number

The Hermitian  $(\mathbf{L}^{(q)} = \overline{\mathbf{L}^{(q)}}^T)$  Magnetic Laplacian encodes edge direction using complex numbers







### Function Name Prediction: <u> Open G</u>raph Benchmark Code 2





#### $\rightarrow$ We outperform previous state of the art by 2.85% (relatively 14.7%)





- I. A new principled **transformer** for **directed graphs** using spectral graph theory
- 2. Specifically, I propose to use the eigenvectors of the **Magnetic Laplacian** as positional encodings for direction and structure awareness.
- Directed graphs can have benefits in domains where sequence representations are common. E.g., for source code, the proposed transformer (a) is invariant w.r.t. meaningless reorderings, and (b) outperforms the prior state of the art by 15%.




 $\rightarrow$  A transformer for directed graphs can handle a vast amount of modalities

Language, time series, ...

Molecules, point clouds, Markov random fields...



Images, PDEs,

Videos, ...



Syntax Trees, Abstract Meaning Representation, Scene Graphs, Knowledge Graphs, Electric Circuits, Logic,





. . .

Outline (Selected Works)

- Transformers Meet Directed Graphs (internship @ Google DeepMind): A Generalization to (Almost) Universal Input Modalities [ICML 2023, DLG-WS @ AAAI 2023]
- 2. Adversarial Robustness in Discrete Domains: Focusing on Graph Neural Networks and Combinatorial Optimization [NeurIPS 2020, 2021, 2022; ICLR 2022, DLG-WS @ AAAI 2021]
- 3. Uncertainty Estimation in Non-standard Settings: Closed-Form Single-Pass Uncertainty Estimation [ICLR 2022 (oral), NeurIPS 2022]



### , Uncertainty in Supervised Learning

#### Input Space



#### **Existing methods**

**X** They often focus on a **single task type**.

- 🗶 They are **overconfident** far from training data.
- \* They require **multiple forward pass** at testing time.
- They require **OOD data at training time**.

### Natural Posterior Network (NatPN)

- I. It applies to many common supervised learning task types.
- 2. It **guarantees** high uncertainty far from training data.
- 3. It requires only a **single foward** pass at testing time.
- 4. It does **not need OOD data** at training time.
- 5. It admits for **posterior predictive distribution in closed-form**

# Exponential Family Distribution

Unified parametrization for Categorical, Normal, Poison, ... distributions based on sufficient statistics  $\chi$ .

Always have a **conjugate prior** which is also an Exponential family distribution.

**Input-dependet update** with number of observations *n* :

 $\mathbf{y}^{(i)} \sim \mathbb{P}(\mathbf{y}^{(i)}|\boldsymbol{\theta}^{(i)})$ 

Target distribution (Aleatoric uncertainty)  $\boldsymbol{\theta}^{(i)} \sim \mathbb{Q}(\boldsymbol{\theta}^{(i)} | \boldsymbol{\chi}^{post,(i)}, \boldsymbol{n}^{post,(i)})$ 

Conjugate-prior distribution (Epistemic uncertainty)  $\chi^{post,(i)} = \frac{n^{prior}\chi^{prior} + n^{(i)}\chi^{(i)}}{n^{prior} + n^{(i)}}$  $n^{post,(i)} = n^{prior} + n^{(i)}$ 







## Examples for Likelihood, Prior and Loss

Likelihood $\mathbb P$	Conjugate Prior $\mathbb{Q}$	Parametrization Mapping $m$	Bayesian Loss (Eq. 5)
$y \sim \operatorname{Cat}(\boldsymbol{p})$	$oldsymbol{p} \sim \mathrm{Dir}(oldsymbol{lpha})$	$egin{array}{lll} oldsymbol{\chi} = oldsymbol{lpha}/n \ n = \sum_c lpha_c \end{array}$	$ \begin{aligned} \mathbf{(i)} &= \psi(\alpha_{y*}^{(i)}) - \psi(\alpha_0^{(i)}) \\ \mathbf{(ii)} &= \log B(\boldsymbol{\alpha}^{(i)}) + (\alpha_0^{(i)} - C)\psi(\alpha_0^{(i)}) - \sum_c (\alpha_c^{(i)} - 1)\psi(\alpha_c^{(i)}) \end{aligned} $
$y \sim \mathcal{N}(\mu, \sigma)$	$\mu, \sigma \sim \mathcal{N}\Gamma^{-1}(\mu_0, \lambda, lpha, eta)$	$oldsymbol{\chi} = egin{pmatrix} \mu_0 \ \mu_0^2 + rac{2eta}{n} \end{pmatrix} \ n = \lambda = 2lpha$	$(\mathbf{i}) = \frac{1}{2} \left( -\frac{\alpha}{\beta} (y - \mu_0)^2 - \frac{1}{\lambda} + \psi(\alpha) - \log\beta - \log 2\pi \right)$ $(\mathbf{i}\mathbf{i}) = \frac{1}{2} + \log\left( (2\pi)^{\frac{1}{2}} \beta^{\frac{3}{2}} \Gamma(\alpha) \right) - \frac{1}{2} \log \lambda + \alpha - (\alpha + \frac{3}{2}) \psi(\alpha)$
$y \sim \operatorname{Poi}(\lambda)$	$\lambda \sim \Gamma(lpha,eta)$	$\chi = lpha / n \ n = eta$	

 $\psi(x)$  and B(x) are the Digamma and Beta function, repsectively



### ه Exemplary Result: Classification



See paper for, e,.g., state-of-the-art OOD detection results



## 💑 Exemplary Result: (Gaussian) Regression









## Thank you for your attention!

