

GNN Explanations that do not Explain and How to find Them

Steve Azzolin[†], STEFANO TESO[†], BRUNO LEPRI[‡],
ANDREA PASSERINI[†], SAGAR MALHOTRA[§]

[†] *University of Trento*

[‡] *Bruno Kessler Foundation*

[§] *TU Wien*

ICLR 2026

Background

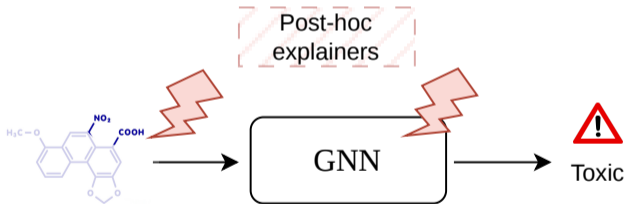
Background

😓 GNNs are black-boxes



Background

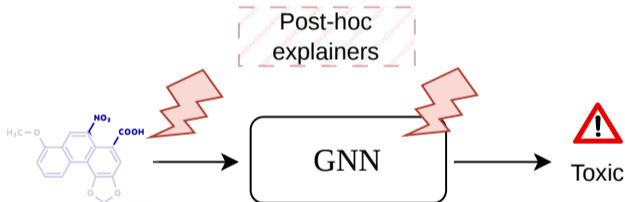
😞 GNNs are black-boxes



Background


😞 GNNs are black-boxes

😞 Post-hoc explainers can be unreliable¹



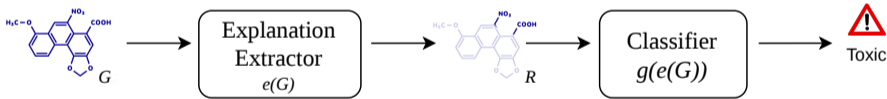
¹C. Rudin. Stop explaining black box machine learning models and use interpretable models instead. Nature Machine Intelligence. 2019

Background: From post-hoc to ante-hoc


 Design GNNs with interpretability in mind

Background: From post-hoc to ante-hoc

- 🤔 Design GNNs with interpretability in mind
- 🤔 **Self-explainable GNNs (SE-GNNs)** generate explanations during inference



Background: From post-hoc to ante-hoc

 How *good* are SE-GNNs' explanations?

Background: From post-hoc to ante-hoc

😬 SE-GNN explanations can be:

- ▶ redundant²

²Tai et al., Redundancy Undermines the Trustworthiness of Self-Interpretable GNNs. ICML 2025

³Azzolin et al., Beyond Topological SE-GNNs: A Formal Explainability Perspective. ICML 2025

⁴Wu et al., Discovering Invariant Rationales for Graph Neural Networks. ICLR 2022

Background: From post-hoc to ante-hoc

😬 SE-GNN explanations can be:

- ▶ redundant²
- ▶ ambiguous³

²Tai et al., Redundancy Undermines the Trustworthiness of Self-Interpretable GNNs. ICML 2025

³Azzolin et al., Beyond Topological SE-GNNs: A Formal Explainability Perspective. ICML 2025

⁴Wu et al., Discovering Invariant Rationales for Graph Neural Networks. ICLR 2022

Background: From post-hoc to ante-hoc

😬 SE-GNN explanations can be:

- ▶ redundant²
- ▶ ambiguous³
- ▶ can be affected by spurious correlations⁴

²Tai et al., Redundancy Undermines the Trustworthiness of Self-Interpretable GNNs. ICML 2025

³Azzolin et al., Beyond Topological SE-GNNs: A Formal Explainability Perspective. ICML 2025

⁴Wu et al., Discovering Invariant Rationales for Graph Neural Networks. ICLR 2022

Background: From post-hoc to ante-hoc

😬 SE-GNN explanations can be:

- ▶ redundant²
- ▶ ambiguous³
- ▶ can be affected by spurious correlations⁴

😬 Are there some pathological failure cases we should be aware of?

Present work: Yes, there are!

²Tai et al., Redundancy Undermines the Trustworthiness of Self-Interpretable GNNs. ICML 2025

³Azzolin et al., Beyond Topological SE-GNNs: A Formal Explainability Perspective. ICML 2025

⁴Wu et al., Discovering Invariant Rationales for Graph Neural Networks. ICLR 2022

Our contribution

- 🙄 SE-GNNs can attain optimal loss by providing **degenerate explanations**
 - ▶ *explanations that do not tell anything about the underlying model*

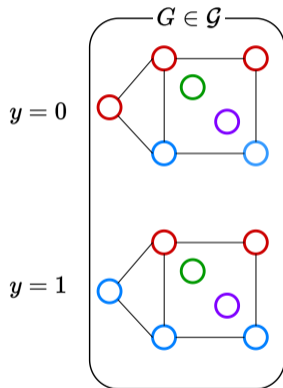
Our contribution

- 🙄 SE-GNNs can attain optimal loss by providing **degenerate explanations**
 - ▶ *explanations that do not tell anything about the underlying model*
- 🙄 Degenerate explanations can go undetected
 - ▶ *faithfulness metrics may not mark them as unfaithful*

SE-GNNs Explanations that do not Explain

Task
 $y = \mathbb{1}\{\#blue > \#red\}$

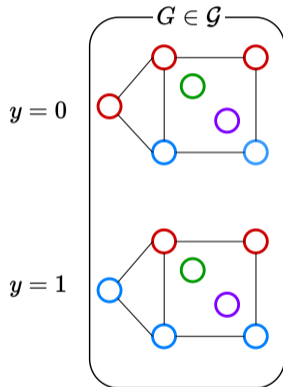
 and 
are uncorrelated with Y



SE-GNNs Explanations that do not Explain

Task
 $y = \mathbb{1}\{\#blue > \#red\}$

○ and ○
are uncorrelated with Y

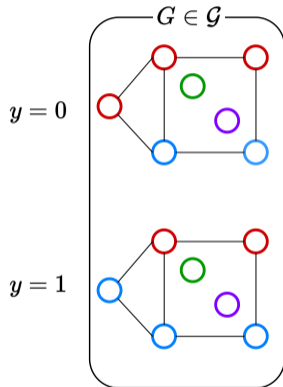


😊 Accurate GNNs must use ○ and ○

SE-GNNs Explanations that do not Explain

Task
 $y = \mathbb{1}\{\#blue > \#red\}$

○ and ○
are uncorrelated with Y



😴 Accurate GNNs must use ○ and ○

😏 Accurate SE-GNNs can provide only ○ and ○ as explanations

This is what we call **degenerate explanations.**

Finding degenerate explanations

? Can we detect when an SE-GNN produces degenerate explanations?

Finding degenerate explanations

? Can we detect when an SE-GNN produces degenerate explanations?

Type	Perturbations \mathcal{I}	Metrics
<i>Nec.</i>	Explanation removal	Fid+, PN
<i>Nec.</i>	Edge removal	RFid+, Nec
<i>Suff.</i>	Complement removal	Fid-, PS, GEF
<i>Suff.</i>	Edge removal	RFid-, SimOAR
<i>Suff.</i>	Complement swap	Suf
<i>Suff.</i>	All	EST (Ours)

Table: Faithfulness metrics.

Finding degenerate explanations

? Can we detect when an SE-GNN produces degenerate explanations?

Dataset	Model	RejRatio _{\mathcal{I}}							
		Fid-	Fid+	Suf	Nec	CF	RFid-	RFid+	EST (ours)
RBGV	SMGNN	12 \pm 19	<u>20</u> \pm 12	00 \pm 00	-	05 \pm 01	00 \pm 00	-	48 \pm 02
	GSAT	12 \pm 21	15 \pm 16	03 \pm 03	-	50 \pm 04	11 \pm 09	-	<u>49</u> \pm 03
	DIR	32 \pm 23	27 \pm 24	03 \pm 03	-	<u>39</u> \pm 08	08 \pm 06	-	48 \pm 03
MNISTsp	SMGNN	88 \pm 03	58 \pm 04	<u>99</u> \pm 01	55 \pm 05	92 \pm 01	<u>99</u> \pm 00	75 \pm 05	100 \pm 00
	GSAT	65 \pm 09	38 \pm 05	<u>99</u> \pm 01	44 \pm 05	95 \pm 02	<u>99</u> \pm 01	61 \pm 04	100 \pm 00
	DIR	93 \pm 03	55 \pm 07	<u>99</u> \pm 01	54 \pm 05	91 \pm 01	<u>99</u> \pm 01	69 \pm 07	100 \pm 00
MUTAG	SMGNN	59 \pm 03	05 \pm 04	99 \pm 00	57 \pm 04	55 \pm 07	72 \pm 03	65 \pm 04	<u>96</u> \pm 01
	GSAT	21 \pm 21	05 \pm 02	99 \pm 00	53 \pm 07	68 \pm 17	70 \pm 14	61 \pm 05	<u>94</u> \pm 05
	DIR	70 \pm 13	04 \pm 02	99 \pm 00	54 \pm 02	69 \pm 08	75 \pm 07	62 \pm 04	<u>97</u> \pm 02
SST2P	SMGNN	08 \pm 02	<u>44</u> \pm 05	14 \pm 05	-	23 \pm 07	04 \pm 01	-	54 \pm 04
	GSAT	<u>51</u> \pm 31	19 \pm 14	07 \pm 10	-	37 \pm 18	14 \pm 03	-	62 \pm 15
	DIR	50 \pm 02	09 \pm 06	12 \pm 06	-	42 \pm 10	05 \pm 01	-	<u>49</u> \pm 03

Table: Previous metrics can fail to reject degenerate explanations.

Conclusions

- ▶ Check out our paper for more details
- ▶ For any question: `steve.azzolin@unitn.it`