

# TopInG: Topologically Interpretable Graph Learning via Persistent Rationale Filtration <sup>1</sup>Cheng Xin\*, <sup>2</sup>Fan Xu\*, <sup>2</sup>Xin Ding, <sup>1</sup>Jie Gao and <sup>2</sup>Jiaxin Ding <sup>1</sup>Rutgers University, <sup>2</sup>Shanghai Jiao Tong University

### **Problem and Contributions**

**Problem**: Learn a model simultaneously predict the label and identify a **rational subgraph**.

- **TopInG** integrates TDA into GNNs for interpreting graphs by learning rationale filtration.
- Propose a new loss, **topological discrepancy**, measuring statistical difference on topological invariants --- persistent homology.
- Provide theoretical guarantees and a tractable approximation of our model.
- Experiments show **TopInG** improves performance on multiple benchmark datasets up to 20%+.

# Challenges

Trade-off between Prediction & Interpretation

Spurious correlations beyond rationale subgraph

[New] Variform Rationale Subgraphs







https://jackal092927.github.io/publication/TopInG\_ICML2025

### Method

Main Idea: modeling an underlying graph generating process determined by an ordering f:G
ightarrow [0,1] consistent with the partition  $G=G_X^*|\;\left|\,G_\epsilon^*
ight|$ i.e.:  $\exists t \approx 0.5, G_{< t} \approx G_X^*, G_{> t} \approx G_{\epsilon}^*$ 

**Persistent Homology:** on a graph filtration  $\mathcal{F}(G) := \{G_{\leq t} \mid t \in f(E)\}$  a persistent homology is a chain of induced homology vector spaces connected by linear maps

 $H_p(\mathcal{F}(G)): 0 \to \cdots \to H_p(G_{\leq t_1}) \to H_p(G_{\leq t_2}) \to H_p(G_{\leq t_3}) \to \cdots \to H_P(G)$ 

Intuitively, using topological discrepancy to enlarge persistent topological gap

 $d_{\text{topo}}(\mathcal{P}(G_X), \mathcal{P}(G_{\epsilon})) \triangleq \inf_{\pi \in \Pi(\mathcal{P}(G_X), \mathcal{P}(G_{\epsilon}))} \mathbb{E}$ 

### Theory

### A tractable lower bound; upper bound by a functional Gromov-Hausdorff distance

 $\max_{\psi\in\Psi} |\mathbb{E}_{P\sim\mathcal{P}(G_X)}[\psi(P)] - \mathbb{E}_{Q\sim\mathcal{P}(G_\epsilon)}[\psi(Q)]| \leq d_{topo}(\mathcal{P}(G_X),\mathcal{P}(G_\epsilon)) \leq C\cdot d_{FGH}(G_X,G_\epsilon)$ 

Theoretical Guarantee: optimized by ground truth rationale subgraphs

**Theorem 3.4.** Assume  $\forall G, |E_X| < |E_{\epsilon}|$ , and  $G_X^*$  is minimal with respect to  $y_G$  in the sense that any subgraph  $G_X \subset G_X^*$  losses some information of label, then  $\hat{\mathcal{L}}(\phi)$  is uniquely optimized by  $f_{\phi}^*(e) = 1\{e \in G_X^*\}$ .



C. Xin and J. Gao acknowledge funding from IIS-22298766, DMS-2220271, DMS-2311064, CRCNS2207440, CCF-2208663 and CCF-2118953. F. Xu and J. Ding were supported by NSF China under Grant No. T2421002, 62202299, 62020106005, 62061146002.

$$\mathbb{E}_{(P,Q)\sim\pi}[d_{\mathrm{B}}(P,Q)]$$

	SingleMotif				MultipleMotif		RealDataset	
Method	<b>BA-2Motifs</b>	<b>BA-HouseGrid</b>	SPmotif0.5	SPMotif0.9	<b>BA-HouseAndGrid</b>	<b>BA-HouseOrGrid</b>	Mutag	Benzene
GNNEXPLAINER	$67.35 \pm 3.29$	$50.73 \pm 0.34$	$62.62 \pm 1.35$	$58.85 \pm 1.93$	$53.04\pm0.38$	$53.21\pm0.36$	$61.98 \pm 5.45$	$48.72\pm0.14$
PGEXPLAINER	$84.59\pm9.09$	$50.92 \pm 1.51$	$69.54 \pm 5.64$	$72.34 \pm 2.91$	$10.36\pm4.37$	$3.14\pm0.01$	$60.91 \pm 17.10$	$4.26\pm0.36$
MATCHEXPLAINER	$86.06 \pm 28.37$	$64.32\pm2.32$	$57.29 \pm 14.35$	$47.29 \pm 13.39$	$81.67\pm0.48$	$79.87 \pm 1.61$	$91.04 \pm 6.59$	$55.65 \pm 1.16$
MAGE	$79.81 \pm 2.27$	$82.69 \pm 4.78$	$76.63\pm0.95$	$74.38\pm0.64$	$99.95\pm0.06$	$99.93\pm0.07$	$99.57\pm0.47$	$96.03\pm0.63$
DIR	$82.78 \pm 10.97$	$65.50 \pm 15.31$	$78.15 \pm 1.32$	$49.08\pm3.66$	$64.96 \pm 14.31$	$59.71 \pm 21.56$	$64.44 \pm 28.81$	$54.08 \pm 13.75$
GSAT	$98.85\pm0.47$	$98.58 \pm 0.59$	$74.49 \pm 4.46$	$65.25 \pm 4.42$	$92.92\pm2.03$	$77.52\pm3.71$	$99.38\pm0.25$	$91.57 \pm 1.48$
GMT-LIN	$97.72\pm0.59$	$85.68 \pm 2.79$	$76.26\pm5.07$	$69.08 \pm 10.14$	$76.12\pm7.47$	$74.36\pm5.41$	$\textbf{99.87} \pm \textbf{0.09}$	$83.90\pm6.07$
TOPING	$\textbf{100.00} \pm \textbf{0.00}$	$\textbf{99.87} \pm \textbf{0.13}$	$\textbf{95.08} \pm \textbf{0.82}$	$\textbf{90.82} \pm \textbf{4.95}$	$\textbf{100.00} \pm \textbf{0.00}$	$\textbf{100.00} \pm \textbf{0.00}$	$96.38 \pm 2.56$	$\textbf{100.00} \pm \textbf{0.00}$

	RealI	Dataset	SpuriousMotif			
	Mutag	Benzene	<b>b=0.5</b>	<b>b=0.7</b>	b=0.9	
DIR	$68.72 \pm 2.51$	$50.67\pm0.93$	$45.49 \pm 3.81$	$41.13 \pm 2.62$	$37.61 \pm 2.02$	
GSAT	$\textbf{98.28} \pm \textbf{0.78}$	$\textbf{100.00} \pm \textbf{0.00}$	$47.45\pm5.87$	$43.57 \pm 2.43$	$45.39\pm5.02$	
<b>GMT-LIN</b>	$91.20 \pm 2.75$	$\textbf{100.00} \pm \textbf{0.00}$	$51.16\pm3.51$	$53.11 \pm 4.12$	$47.60\pm2.06$	
TOPING	$92.92 \pm 7.02$	$\textbf{100.00} \pm \textbf{0.00}$	$\textbf{79.30} \pm \textbf{3.92}$	$\textbf{75.46} \pm \textbf{7.62}$	$\textbf{65.64} \pm \textbf{4.98}$	





# Experiments

Table 1: Interpretation Performance (AUC) on benchmark datasets.

Table 2: Prediction Accuracy (Acc) on benchmark datasets.

### Visualization

$\checkmark$		~~ V	P	8
H A A	A	A A A		
H H	HAR AND	H H	A	HAN A