

Perturb More, Trap More: Understanding Behaviors of Graph Neural Networks

Chaojie Ji^{a,b,*}, Ruxin Wang^{a,b,*}, Hongyan Wu^{a,b,**}

^a *Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences, Shenzhen, China*

^b *Joint Engineering Research Center for Health Big Data Intelligent Analysis Technology*

Abstract

While graph neural networks (GNNs) have shown a great potential in various tasks on graph, the lack of transparency has hindered understanding how GNNs arrived at its predictions. Although few explainers for GNNs are explored, the consideration of local fidelity, indicating how the model behaves around an instance should be predicted, is neglected. In this paper, we first propose a novel post-hoc framework based on local fidelity for any trained GNNs - **TraP2**, which can generate a high-fidelity explanation. Considering that both relevant graph structure and important features inside each node need to be highlighted, a three-layer architecture in TraP2 is designed: i) interpretation domain are defined by **Translation** layer in advance; ii) local predictive behavior of GNNs being explained are probed and monitored by **Perturbation** layer, in which multiple perturbations for graph structure and feature-level are conducted in interpretation domain; iii) high faithful explanations are generated by fitting the local decision boundary through **Paraphrase** layer. Finally, TraP2 is evaluated on six benchmark datasets based on five desired attributions: accuracy, fidelity, decisiveness, insight and inspiration, which achieves 10.2% higher explanation accuracy than the state-of-the-art methods.

*Equal contribution.

**Corresponding authors at: Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences, Shenzhen, China and Joint Engineering Research Center for Health Big Data Intelligent Analysis Technology.

Email addresses: cj.ji@siat.ac.cn (Chaojie Ji), rx.wang@siat.ac.cn (Ruxin Wang), hy.wu@siat.ac.cn (Hongyan Wu)

Keywords:

Graph Neural Network, Explainability

1. Introduction

The development of deep neural networks has shown outstanding performances in many domains [1, 2, 3], which is partially attributed to the effectiveness of mining latent representations from Euclidean domain. By using regular grid convolution, convolutional neural networks (CNNs) effectively extract important semantic information from Euclidean data. However, there is an increasing number of applications which data are represented as graphs. Recently, Graph neural networks (GNNs) [4, 5, 6, 7] have been proposed and achieved breakthroughs in processing various graph structure data, such as point clouds, social network, chemistry molecules solved some classical field problems such as node classification [8] and graph classification [9]. With the advent of more high-precision models, requirements for a better understanding of the inner workings of a model are getting more attention. It is important to know the reason why this model produces such predictions. From the aspect of application, models need to provide transparent and high-precision solutions especially for some key scenarios, such as security, economy and healthcare. Meanwhile, interpretable approach can provide insight for improving models, and show underlying rules that are overlooked. Therefore, higher demands on the performances are raised for explainer: 1) the explainers need provide accurate, reliable and stable explanation for original model. 2) excellent interpretable approaches are expected to provide insight for model cognition and reveal potential rules that are neglected. However, the design of end-to-end network architecture reduces the transparency of the model, which hinders people to comprehend them.

To improve the transparency of deep neural networks, many researches of explanation technologies have been introduced and applied in recent years [10, 11, 12, 13]. Regrettably, existing explanation methods encountered barriers when they are applied to GNNs due to the particularity of graph structure. In recent two years, few excellent explanation approaches for GNNs have been studied. A qualified interpretation model should provide accurate explanation and be faithful to original model. Although learning a completely faithful explanation is usually impossible, a meaningful interpretation in the vicinity of prediction being explained is possible [14]. Unfortunately, taking

node classification as an example, prior researches only use the GNN’s prediction behavior information of node being explained to generate explanations, which neglects the other local behaviors near the decision boundary of the node. These behavior informations are very useful for providing a local faithful explanation. To our best knowledge, there is no method considering the “vicinity” in graph-structured data for a post-hoc explanation based on local fidelity at present. In addition, each feature component is node-dependent that each feature can make different contribution for distinct nodes and the prediction behaviors of GNNs. It is crucial for evaluating and integrating them for an accurate and reliable explanation. However, existing methods ignore this. For example, the GNNExplainer [15] only coarsely provides a shared score for the feature components of all nodes, which leads to a sub-optimal interpretation effect.

To address above problems, we propose a novel post-hoc and model-agnostic explanation framework for GNNs, which provides a “broad” insight of original model (Figure 1). To achieve this, the proposed approach includes a three-layer architecture design which is named **TraP2**: i) **T**ranslation Layer is adopted to realize the transformation from the original problem to interpretation domain according to different tasks. ii) To “trap” richer predictive behaviors near a local decision boundary of the object being explained, the **P**erturbation Layer probes and monitors local behaviors by specially designed strategy from both graph structure-level and feature-level. Meanwhile, a novel perturbation energy level is proposed for measuring the obtained attention of each perturbation. iii) Based on perturbed instances, **P**araphrase Layer finds high-faithful explanations by fitting the local decision boundary, which also provides insight in both graph structure-level and node-dependent feature-level. Finally, to verify the effectiveness of the presented model, we evaluate it on multiple datasets [15] for node and graph classification, and compare our results with other state-of-the-art methods. The results demonstrate the effectiveness of our work from accuracy, fidelity and contrastivity, which outperforms other approaches for all tasks. In summary, the main contributions of this paper can be summarized as follows:

- We first put forward a novel post-hoc framework **TraP2** based on local fidelity of any GNN models for different recognition tasks, which generates high-fidelity explanations.
- Novel perturbation strategy and perturbation effect estimation method designed for graph data are proposed. Furthermore, our model provides

a more fine-grained explanation on node-dependent feature-level than prior works.

- Compared with the state-of-the-art GNN explanation approaches, the proposed method achieves the top performance on multiple benchmark datasets.

The remainder of this paper is organized as follows: Section 2 is the related works containing the introduction of the graph neural networks, current non-graph and graph neural networks-based interpretability methods. In Section 3, we introduce the details of our proposed explanation approach. Experimental results for node and graph classification tasks are given in Section 4. Finally, Section 5 concludes this paper and prospects some future works.

2. Related Work

2.1. Graph neural networks

In recent years, graph neural networks have been successfully applied to a wide variety of fields such as computer vision [16, 17], natural language processing [18], recommender systems [19] and healthcare [20]. GNNs effectively handle the complex relationship between objects in the graph structure following a neighborhood aggregation scheme, where the features of a node is computed by recursively assembling and transforming features of its local neighbors [21, 22]. They can be divided into two categories: spectral [23, 24] and spatial methods [25, 26, 27]. The spectral method is defined via graph Fourier transform and convolution theorem, which takes the graph Laplace matrix as an important tool and does not explicitly use the information propagation mechanism on the graph. Bruna et al. [28] proposed a method to conduct convolution in the spectral domain adopting the Fourier basis of a given graph. Levie et al. [29] introduced a new spectral domain convolutional architecture using a new class of parametric rational complex functions (Cayley polynomials) that can specialize on frequency bands of interest. On the other hand, spatial approaches define convolution operations in the vertex domain, operating on spatially local neighborhood nodes. For instance, Ying et al. [30] proposed a differentiable graph pooling method for GNNs, which can be used to obtain a representation of an entire graph by summing the features of all nodes in the graph. Veličković et al. [31] presented graph

attention networks operated on graph-structured data, assigning different weights to different nodes within a neighborhood.

2.2. Non-graph neural networks interpretability methods

Many well-studied explanation techniques employ gradient-based back-propagation to calculate saliency maps for original input. Some prominent methods of this category include Class Activation Mapping (CAM) [32], Gradient-weighted Class Activation Mapping (Grad-CAM) [33], (Grad-CAM++) [34], Excitation Back-Propagation (EB) [35] and Layer-wise Relevance Propagation (LRP) [36]. The CAM and its generalization Grad-CAM measure the linear combination of each layer’s activations and class-specific weights or gradients of original model. The EB improves gradient maps by introducing contrastive top-down attention. And the LRP adopts layer-wise relevance propagation to achieve a pixel-wise decomposition. The above approaches introduce different backpropagation heuristics, which can focus on salient notions of input data. However, they are not model-agnostic, with most of them being limited to original network framework and/or many necessary modifications of model structure [37]. LIME [14] is a representative system of model-agnostic approaches, which adopts the self-explanatory linear regression to local area and pinpoints important features based on the regression coefficients. The SHAP model [38] uses the Shapley values of a conditional expectation function of the original model to measure the importance of each feature for a particular prediction.

2.3. Graph neural networks interpretability methods

To the best of our knowledge, few explainers for GNNs are explored recently. Pope et al. [39] extended the gradient-based saliency map methods to GNNs, which utilizes the network parameters and classifier output of original GNNs to construct the activation response of the corresponding neurons. Baldassarre et al. [40] employed two main classes of techniques, gradient-based and decomposition-based, to learn important components of input that also relies on propagating gradients/relevance from the output to the input of original model. Ying et al. [15] proposed an explanation method by identifying a small subgraph structure and a subset of node features that maximizes the mutual information with original GNNs prediction in entire input graph. Although the above works have made breakthrough progress in GNNs interpretation, the consideration of local fidelity and fine-grained

$\hat{V} \subseteq V$ and $\hat{E} \subseteq E$ highlights the explanation graph structure. Furthermore, the interpretable features of nodes inside \hat{V} are marked as \hat{X} .

For simplifying the following formulation, we explicitly define two functions: generating reachability matrix and fetching specified elements from an adjacency matrix. Reachability matrix records the connectivity between nodes within k -hop. It is able to be derived from multiplying the adjacency matrix A by itself. We formulate the function of generating it as $[A^k]$. Fetching an element of the i th row and j th column from a given matrix M is formulated as $[M]_{i,j}$. In particular, $[M]_{i,:}$ represents fetching i th row in M .

3.2. Graph Neural Networks

GNNs update the representation of each node through summarizing the local information similar to convolution operation in CNNs. At each layer l , graph structure (V, E) is kept unchanged, only the features of the graph are updated:

$$m_i^{l+1} = \sum_{j \in \mathcal{N}_i} \mathcal{F}(h_i^l, h_j^l, e_{ij}), \quad h_i^{l+1} = \mathcal{G}(m_i^{l+1}, h_i^l) \quad (1)$$

where \mathcal{N}_i defines the neighbor nodes around node i . \mathcal{F} and \mathcal{G} are message and update function respectively. The representation of node i in layer l is denoted as h_i^l . h_i^0 is initialized as x_i . e_{ij} symbolizes the type of edge from node i to j . The updated representation after final layer L can be mapped as $f : h^L \Rightarrow Y$ for a specific task, i.e. node classification (regression), graph classification (regression) and etc.. Formally, GNNs can be formulated as $f(A, X)$.

4. Method

In this section, we describe the uniform framework of our model - TraP2. TraP2 has three components: 1) Translation Layer: accessible subgraphs to be explained are remained from complete input graphs; 2) Perturbation Layer: translated subgraphs are respectively perturbed in aspects of graph structure and node feature. Meanwhile, the degree of these perturbation is assessed as ‘‘perturbation energy level’’; 3) Paraphrase Layer: a local faithful interpretation is trained to explain the behaviors of GNNs, which is shown in Figure 1.

In this section, we first introduce the explanation process of node classification on node i and then extend it into other graph tasks.

4.1. Translation Layer: Transference from Source Domain to Interpretation Domain

In some cases, the domain to be explained for one node is not the entire graph (Source Domain) as shown in Figure 1 (a). To be more precise, GNNs with k layers can only aggregate messages from k -hop, e.g. 2-3, neighbors [5]. According to Equation (1), final representation of node i , h_i^l , is obtained in l recursive updates. In each update, nodes located in one-hop farther from node i are aggregated. Thus only nodes within l -hop from node i is considered in GNNs. It directly results in that the feasible region for explanation is merely a small subgraph (Interpretation Domain) as shown in Figure 1 (b) instead of entire graph Figure 1 (a). Benefited from it, the solution space can be potentially shrank.

Inspired by it, translation layer is applied to transforms the source graph-structured domain into a limited interpretable domain for node i :

$$V^{\mathcal{I}} = \{x | x > 0, \forall x \in [A^k]_{i,:}\} \quad (2)$$

$$A^{\mathcal{I}} = [A]_{V^{\mathcal{I}} \times V^{\mathcal{I}}} \in \{0, 1\}^{\hat{n} \times \hat{n}} \quad (3)$$

where $V^{\mathcal{I}}$ and $A^{\mathcal{I}}$ denote the node set and corresponding adjacency matrix in interpretation domain respectively. \hat{n} is the number of elements in set $V^{\mathcal{I}}$.

4.2. Perturbation Layer: Turbulence in Interpretation Domain

In order to realize a local-fidelity based explanation, the strategies of perturbation are proposed to probe the responses of GNNs in local vicinity. Accordingly, a series of novel perturbation patterns are designed for graph. These perturbed instances can be distinguished from the aspects of both graph structure and feature as shown in Figure 1 (c). In particular, we introduce the concept of perturbation energy level that magnitude of each disturbance is quantified. And it'll be further delivered into paraphrase layer for establishing attention for each perturbation. Finally we monitor and record the corresponding behavior response of GNNs to be explained in the disturbances.

4.2.1. Perturbation on Graph Structure

We first define an action set on graph structure: i) adding new edges between nodes in $V^{\mathcal{I}}$, ii) removing existing edges from $A^{\mathcal{I}}$.

A random action variable, z^E , is applied to trigger these two actions alternatively:

$$z^E \sim \text{Bernoulli}(1, p_1) \quad (4)$$

where p_1 is the probability for actions. As a graph can be represented as a binary adjacency matrix which can be deformed through alternating 0 and 1 for each element in it. The composition of a perturbed graph A^P under different patterns can be decomposed as:

$$[A^P]_{j,k} = \begin{cases} [A^I]_{j,k} \oplus (z^E - 1), & \text{Adding Pattern} \\ [A^I]_{j,k} z^E, & \text{Removing Pattern} \\ [A^I]_{j,k} \oplus z^E, & \text{Adding \& Removing} \end{cases} \quad (5)$$

in which $[A^P]_{j,k} \in \{0, 1\}$ and \oplus denotes *xor* operation. Actions of adding and removing edges are alternatively or simultaneously regarded in these three patterns respectively.

Furthermore, specific constraint on perturbation patterns can be additionally imposed. For an instance, if A^I is relatively sparse that number of nodes located within 1-hop from node i is extremely smaller than farther hops, an arbitrary perturbation on A^I probably results in a great impact - far from the vicinity. Given a specified case that all directly connected node around node i is a single vertex j which is further expended with many other farther nodes, it implies that node j provide a most valuable clue for analyzing node i due to the smallest hop distance. In such situation, large indiscriminate perturbations probably result in constant absences of node j and thus explainer is forced to neglect this informative node in paraphrase layer. To solve it, the perturbations occurred on the edges that directly connect with node i should be prohibited as $[A^P]_{i,j} = [A^I]_{i,j}, \forall (i, j) \in E$.

4.2.2. Perturbation on Feature

Similarly, we design a perturbation pattern for feature by scaling or masking its representation. z^F is applied as the random action variable for perturbation on feature:

$$z^F \sim \begin{cases} \text{Bernoulli}(1, p_2), & \text{Masking Pattern} \\ \mathcal{N}(0, 1), & \text{Scaling Pattern} \end{cases} \quad (6)$$

Features of node i and z^F are combined as:

$$[X^P]_{i,d} = [X]_{i,d} z^F \quad (7)$$

4.2.3. Perturbation Energy Level

“Perturbation energy level” explicitly quantifies the energy consumption of perturbation according to both graph structures and features. Obviously, more complicated perturbations consume a larger amount of energy, which always introduces more severe deformation than more edges are removed and added and features inside nodes are masked. Then corresponding response from GNNs also quite differs from the original graph. In contrast, samples produced by slighter perturbations with smaller energy consumption always locate in an immediate vicinity of the original samples. As a result, observation on smaller perturbations with slightly altered prediction provides a more important clue to track the “logic” of GNNs, emphasized in learning stage of paraphrase layer.

For the respect of graph structure, we assume that the combination of the distance between perturbed position and node i , and the deformation degree - scale of altered edges - jointly indicate the energy level. Concretely, the closer perturbation and larger deformation usually consume more energy and vice versa.

We define the measurement of perturbation distance as hop value and formulate a normalized coefficient for k -hop as:

$$w_k = \frac{K}{k + 1} \quad (8)$$

$$\alpha_k = \frac{e^{(w_k)}}{\sum_{i \in \{1, \dots, K\}} e^{(w_i)}} \quad (9)$$

where K is the pre-defined maximum, e.g. number of layers in original GNNs. Finally, we obtain the energy level by combining the coefficients with the deformation:

$$\gamma_A = \sum_{k=1}^K \alpha_k \text{sim}([(A^P)^k]_{i,:}, [(A^I)^k]_{i,:}) \quad (10)$$

For the energy level of feature, we measure the similarity between original and perturbed node feature:

$$\gamma_X = \sum_{i=1}^N \text{sim}([X^P]_{i,:}, [X]_{i,:}) \quad (11)$$

where

$$\text{sim}(u, v) = e^{\left(\frac{-\cosine(u, v)^2}{\delta^2}\right)} \quad (12)$$

in which δ is the width with distance function $\text{cosine}(\cdot)$. It suggests that higher energy is consumed with more difference in node features.

To this end, complete energy level is defined as follows:

$$\gamma = \lambda_A \gamma_A + \lambda_X \gamma_X \quad (13)$$

where λ_A and λ_X control the balance between perturbation on graph structure and feature.

4.2.4. Monitor on Multiple Perturbations

Given a single instance, the complex decision boundaries of GNNs being explained can be hardly identified by limited witnesses. Thus multiple perturbations are applied and we formulate a series of independent instances perturbed as $\{(A_{(j)}^{\mathcal{P}}, X_{(j)}^{\mathcal{P}}, \gamma_{(j)})\}_{j=1}^m$ in which m is the frequency of the perturbations. Meanwhile, for each perturbation, we constantly monitor and record the behavior feedback $f(y|A_{(*)}^{\mathcal{P}}, X_{(*)}^{\mathcal{P}})$ of the GNNs being explained. $*$ represents a perturbation.

4.3. Paraphrase Layer: Explanation on Interpretation Domain

Once multiple behavior responses of GNNs under various perturbations are collected and learnt sequentially, a local decision boundary can be identified by the explainer.

4.3.1. Learning Phase

In order to learn a local faithful explanation, we denote an explainer as g which can be any kind of potentially explainable models, including linear models and non-linear models. $w_\theta \in \mathbb{R}^{\hat{n}d}$ is the trainable parameters, the size of which is completely correlated with the scale of $A_{(*)}^{\mathcal{I}}$ and X . d is unchangeable and derived from the number of features inside X . As mentioned in section 4.1, the size of $A_{(*)}^{\mathcal{P}}$ is largely reduced from A according to a transformation executed by the translation layer. That is to say that \hat{n} equals to the number of vertices within limited hop - defined by the GNNs to be explained - from the explained node. Accordingly, it ensures that $\hat{n} \ll n$. And g has ability to generate explanations for an even relatively large graph.

The explainer g is calculated as:

$$g(y|A_{(*)}^{\mathcal{P}}, X_{(*)}^{\mathcal{P}}; w_\theta) = \sigma(w_\theta(\|_{j=1}^n([(A_{(*)}^{\mathcal{P}})^k]_{i,j} \cdot [X_{(*)}^{\mathcal{P}}]_{j,:}))) \quad (14)$$

where symbol $\|$ indicates a concatenation operation, \cdot is a scalar multiplication and σ denotes a nonlinearity function.

In addition, unlike GNNExplainer, the explanation of TraP2 on feature is node-dependent as shown in Figure 1 (d).

We combine explainer g with GNNs being explained f to fit a local decision boundary:

$$\min_{w_\theta} \sum_{j=1}^m \sum_{c=1}^C \frac{1}{\gamma_{(j)}} \mathcal{L}(f(y = c|A_{(j)}^{\mathcal{P}}, X_{(j)}^{\mathcal{P}}), g(y = c|A_{(j)}^{\mathcal{P}}, X_{(j)}^{\mathcal{P}})) + \lambda\Omega(w_\theta) \quad (15)$$

where \mathcal{L} is a measurement of the difference between explainer and GNNs in the locality. Fewer energy consumption $\gamma_{(j)}$ causes more attention on this perturbation in that they are closer to the local decision boundary. $\Omega(w_\theta)$ is regarded as a regularization term to encourage w_θ to be discrete and be interpretable by humans.

4.3.2. Identifying Phase

In our work, the obtained parameters w_θ directly represent the contribution score of a single feature for a node. Further we compute the contribution (explanation) of node j with w_θ :

$$I_j^i = \sum_{q=j*d}^{(j+1)*d} |[w_\theta]_q| \quad (16)$$

where $|[w_\theta]_q|$ indicates the contribution, an absolute value, from the d th feature inside node j . A main goal that explainers aim at is locating nodes that relatively maximum extents of contributions are made. Therefore, the direction of contributions, either positive or negative, have no effect on quantifying the extent.

At last, by sorting of the contribution scores across $\{I_1^i, \dots, I_n^i\}$ and $\{|[w_\theta]_1|, \dots, |[w_\theta]_{nd}|\}$, we can achieve the $\hat{V} \subseteq V^{\mathcal{I}}$ and \hat{X} with higher scores respectively. Furthermore, \hat{E} can be extracted from existing edges inside E and nodes belonging to \hat{V} . Thus, a complete explanation of node i , $\hat{G} = (\hat{V}, \hat{E}, \hat{X})$, is determined.

4.4. Extension on Graph Classification Task

Our proposed TraP2 can not only explain on node classification but also other graph machine learning tasks, e.g. graph classification.

In node tasks, only one node i need to be explained, thus TraP2 appropriately executes one time. However, every node in the graph all makes contributions to graph prediction. According to it, for each node, we independently perform a translation, perturbation and paraphrase as mentioned above. Specifically, there are two small changes in paraphrase layer. i) The c in Equation (15) belongs to label of graph instead of node; ii) In identifying stage, the contribution scores of each node j must be pooled across all nodes:

$$I_j = \frac{1}{n} \sum_{i=1}^n I_j^i \quad (17)$$

5. Experiment

In this section, we conduct experiments on two kinds of tasks, node and graph classification to evaluate the performance of TraP2.

5.1. Datasets

For these tasks, we follow existing study to apply the same benchmark datasets [15].

- **BA-SHAPES** is a Barabási-Albert (BA) graph with 300 nodes and 80 five-node house attachments which are randomly attached. The classes are determined by that nodes locate on the top, middle, bottom or out of a house.
- **BA-COMMUNITY** consists of two BA graphs as in Figure 2. The definition of class for each graph is consistent with BA-SHAPES. In addition, normally distributed features are assigned for each node.
- **TREE-CYCLE** is a 8-level balanced binary tree randomly attached with 80 six-node cycles. The classes are distinguished by whether nodes locate on the cycles.
- **TREE-GRID** is similar with TREE-CYCLE except of that six-node cycles are replaced by 3-by-3 grids.
- **MUTAG** is a dataset composed of 4,337 molecule graphs which are labeled with the mutagenic effect on the Gram-negative bacterium *S.typhimurium* [41].

- **REDDIT-BINARY** is a dataset of 2,000 graphs that online discussion threads on Reddit are recorded. Inside each graph, nodes indicate users and edges stand for the reply between users. The labels of the graphs are the types of user interactions [42].

5.2. Baseline Methods

We compare TraP2 with four baseline methods:

- **Random** is an approach that all explanation elements are randomly generated.
- **Greedy** is a method that the nodes causing the highest accuracy difference from original GNNs are added iteratively, when the edges between these nodes and the node being explained are masked one by one.
- **Grad** is similar to a saliency map method that the gradient is derived from the loss of GNNs being explained [39].
- **GNNE explainer** explores a small subgraph of the entire graph to maximize the mutual information with the prediction of original GNNs.

5.3. Implementation

The implementation and hyperparameters¹ for GNNs to be explained and GNNE explainer are completely derived from the original work. For TraP2, in translation layer, we set k in Equation (2) to 3 corresponding to the trained GNNs. Sample rate p_1 and p_2 in Equation (4) and (6) are respectively set as 0.5 and 0.8 for all datasets. The width δ with distance function in Equation (12) is assigned as 25. We choose 1 for both λ_A and λ_X in Equation (13). The frequency of perturbation is 1500. L1 regularization is selected as our regularization term $\Omega(w_\theta)$ in Equation (15). In paraphrase layer, g is trained in 300 epochs with 0.01 learning rate. To account for the property of TREE-CYCLE and TREE-GRID datasets that number of nodes within 1-hop is smaller than 2 and 3-hops, we constrict the perturbation on graph structure that edges belonging to 1-hop from the node being explained are kept unchanged as mentioned in Section 4.2.

¹<https://github.com/RexYing/gnn-model-explainer>

5.4. Investigation on TraP2

In our work, five desired attributions: accuracy, fidelity, decisiveness, insight and inspiration are investigated for evaluating the performance of all explainers.

Dataset	Random	Greedy	Grad	GNNExplainer	TraP2
BA-SHAPES	1.9	60.8	81.4	81.7	81.9
BA-COMMUNITY	1.6	65.3	66.4	68.6	71.2
TREE-CYCLE	47.7	56.1	71.7	71.8	72.2
TREE-GRID	74.6	75.5	70.2	76.3	86.5

Table 1: Explanation Accuracy on four benchmark datasets.

Dataset	Metrics	Grad	GNNExplainer	TraP2
BA-SHAPES	Fidelity	0.042	0.030	0.020
	Contrast	1.01	1.47	6.98
BA-COMMUNITY	Fidelity	0.15	0.11	0.09
	Contrast	1.06	1.51	5.08
TREE-CYCLE	Fidelity	0.275	0.406	0.272
	Contrast	1.00	1.06	3.51
TREE-GRID	Fidelity	0.718	0.069	0.057
	Contrast	1.01	1.45	4.50

Table 2: Explanation Fidelity and Contrast on four benchmark datasets.

5.4.1. Question 1: Can TraP2 make an accurate explanation for ground-truth knowledge?

One of the most essential criterion of an explanation is the accuracy that the interpretation should match ground truth knowledge exactly. We define the ground truth explanation for BA-SHAPES, BA-COMMUNITY, TREE-CYCLE and TREE-GRID as their attachments, i.e. house, cycle and grid shape. For each node being explained, we train an explainer and sort the contribution scores of each node. We only preserve the top n nodes - number of the nodes inside ground truth - with higher scores and calculate the matching accuracy rate with the ground truth. Intuitively, a better explanation is expected to achieve a higher accuracy. The results of node classification on four datasets are shown in Table 1. TraP2 outperforms

other baseline methods on all datasets, even up to 10.2% higher than the state-of-the-art approach in TREE-GRID that most complicated structures of grid and tree are combined tightly. It is noteworthy that Random method surprisingly obtains a reasonable accuracy in TREE-CYCLE and TREE-GRID. The reason is that the scale of interpretation domain is relatively small.

5.4.2. Question 2: Is TraP2 faithful to the authentic prediction?

A faithful explainer is attributed as a repeater that the original behaviors of GNNs can be reproduced without loss. That is to say that either entire graph or explanation subgraph is fed into trained GNNs, the outputs ought to be identical. However, there is no guarantee that ground truth explanations are consistency with the internal logic of GNNs explained. From another aspect, fed a subgraph composed of ground truth explanations, GNNs may produce predictions which is extremely different from original outputs. In a conclusion, explainers with high accuracy are not the equivalent of the interpreters with high fidelity.

We formally formulate the fidelity as the absolute difference between the prediction from the entire graph and the explanation subgraph consisting of the n nodes from Question 1. Smaller difference explicitly indicates the more faithful explanation. Extensive experiments are conducted on TraP2 and two state-of-the-art approaches with relatively high accuracy.

The fidelity scores are listed in Table 2. For a deeper analysis, explainers with high accuracy are not the equivalent of the interpreters with high fidelity. Although GNNExplainer exceeds Grad on all datasets under accuracy metric, it does not outperform Grad on fidelity for all cases. In contrast, Grad with lowest accuracy unexpectedly defeats the GNNExplainer in TREE-CYCLE dataset. For our method, TraP2 surpasses all alternative methods in four datasets, which proves the highest fidelity of TraP2 among them.

5.4.3. Question 3: Does TraP2 decisively provide explanations?

Another valuable characteristic of an explainer is whether it can make an explanation without hesitation. As a decisive explainer, it is expected to have an ability to firmly determine a distinguished boundary among the explanation elements with high-contrast contribution value.

We define the contrastivity by subtracting the average contribution score of the nodes outside of the n nodes, extracted from Question 1, from the

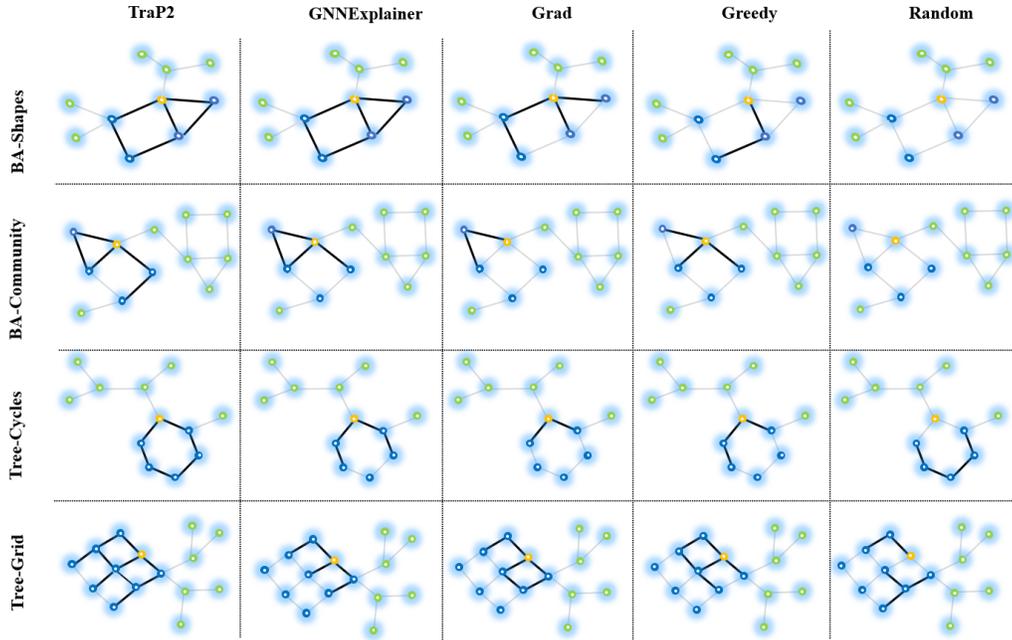


Figure 2: Visual comparisons of the node classification examples on four datasets. The blue nodes indicate the ground truth, the interpreted node is marked in orange. Best viewed on a computer screen.

lowest score inside of the n nodes. A higher contrastivity implies that a clear boundary is simple to identify.

As shown in Table 2, TraP2 achieves the highest contrastivity value across all datasets and approaches. Especially, the superiority is apparently clear that the improvements on GNNExplainer and Grad are about 3.6 and 4.9 times on average respectively. Compared with the low contrastivity of Grad and GNNExplainer, TraP2 can be regarded as a definitely high-contrast method among them.

5.4.4. Question 4: Can TraP2 make insightful explanation for models' predictions?

In addition to accurately identifying the relevant subgraph structure, an explainer should highlight the most meaningful features inside the node and provide a more insightful explanation for understanding model. We compare the selected important features in each node for GNN's prediction with different approaches. In BA-Community dataset (Figure 3), the TraP2 cor-

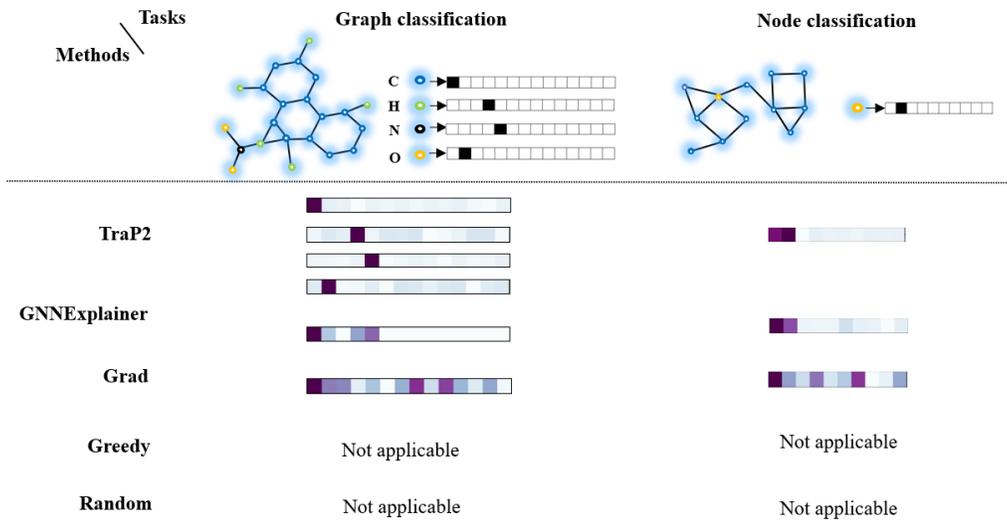


Figure 3: Visual comparisons of the node features on MUTAG and BA-Community datasets. Best viewed on a computer screen.

rectly recognizes the important feature component in related nodes. Furthermore, a toy experiment is conducted using same mechanism in GNNExplainer for node-independent feature contribution, which achieves the accuracy of 70.3 and still higher than other baselines. Compared to TraP2 with node-dependent feature contribution, accuracy decreases by 0.9%. It implies that the joint optimization of subgraph structure and node-dependent features in our method is helpful to improve the accuracy of explanation. Similarly, in Figure 3, the atoms (i.e. C, H, O and N) of molecule have different features in graph. These node features of each atom are also identified by TraP2 accurately. In comparison, GNNExplainer also captures the most important four positions of features, but does not achieve fine-grained discrimination. However, other methods can hardly identify or give incorrect explanations inside the node.

5.4.5. Question 5: How TraP2 find inspirational pattern from graphs?

Supposing we have no prior knowledge about the GNNs - treating the original model as a black box, can the explainers provide us with a meaningful discovery? Obviously Grad requires exposed internal details of explained GNNs, while TraP2 and other baselines can be all categorized as model-agnostic approaches that no details are accessible for explainers. To answer

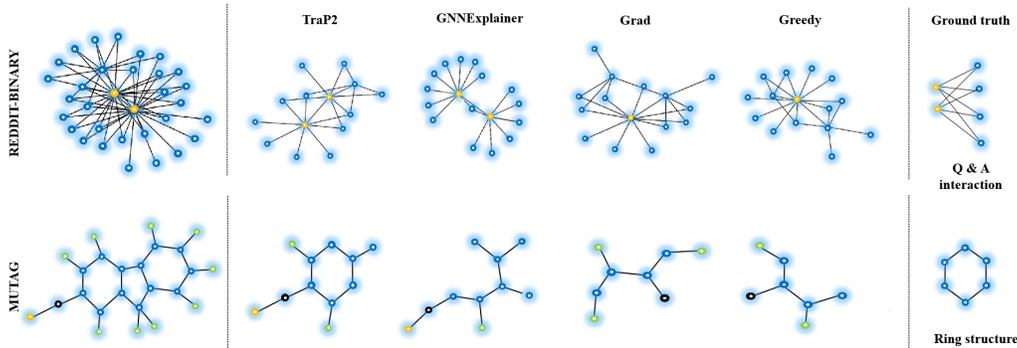


Figure 4: Visual comparisons of the graph classification examples on two datasets. Best viewed on a computer screen.

this question, we revisit the visualization results. For node classification, Figure 2 shows the different methods for four different datasets and highlights the explanation subgraphs of each method. From the figure, it can be seen that TraP2 achieves the best identification performances of key structures of house, cycle and grid. Specifically, Grad and Greedy strategy can recognize part real structures, but they are unsatisfactory because they can not provide an intuitive and understandable subgraph structure. Compared with these two models, GNNExplainer locates more crucial nodes. However, it still encounters some incomplete solutions. But the prediction effect of the random method varies greatly each time, which is reflected in both the accuracy and unstable subgraph structure. As illustrated in Figure 4, for REDDIT-BINARY dataset, the task is to identify whether a given graph belongs to a question/answer-based (Q & A) community. TraP2 automatically finds two dense interaction patterns from complex relationship network. The original network actually represents a Q & A interaction between two experts (the yellow nodes) and multiple visitors (the blue nodes). Similarly, the carbon ring is correctly identified by TraP2, which indicates the mutagenic factor in MUTAG dataset.

6. Conclusion

In this paper, based on local fidelity, we propose a novel explanation framework TraP2, in which local behaviors probed from perturbed instances in vicinity are trapped to generate a high-faithful explanation. Our explanation further explain each task on not only node but also features inside each

nodes. We also propose the fidelity and the contrastivity evaluation metrics to validate the explanation performances. Extensive comparative evaluations on multiple datasets are implemented, which validates the superiority of TraP2 over several state-of-the-art explainers for GNNs. Overall, TraP2 is well adapted and applied to different interpretation tasks, which provides better explanation performance for GNNs. Based on the outstanding performance of our work, we will extend our TraP2 to support more graph mining task such as link prediction and graph generation.

References

- [1] V. Badrinarayanan, A. Kendall, R. Cipolla, Segnet: A deep convolutional encoder-decoder architecture for image segmentation, *IEEE Trans. Pattern Anal. Mach. Intell.* 39 (12) (2017) 2481–2495.
- [2] Z. Zhao, P. Zheng, S. Xu, X. Wu, Object detection with deep learning: A review, *IEEE Trans. Neural Netw. Learn. Syst.* 30 (11) (2019) 3212–3232.
- [3] B. Alshemali, J. Kalita, Improving the reliability of deep neural networks in nlp: A review, *Knowledge-Based Syst.* 191 (2020) 105210.
- [4] W. Hamilton, Z. Ying, J. Leskovec, Inductive representation learning on large graphs, in: *Proceedings of the Advances in Neural Information Processing Systems*, 2017, pp. 1024–1034.
- [5] T. N. Kipf, M. Welling, Semi-supervised classification with graph convolutional networks, *arXiv preprint arXiv:1609.02907* (2016).
- [6] I. Chami, Z. Ying, C. Ré, J. Leskovec, Hyperbolic graph convolutional neural networks, in: *Proceedings of the Advances in Neural Information Processing Systems*, 2019, pp. 4869–4880.
- [7] C. Ji, R. Wang, R. Zhu, Y. Cai, H. Wu, Hopgat: Hop-aware supervision graph attention networks for sparsely labeled graphs, *arXiv preprint arXiv:2004.04333* (2020).
- [8] S. Abu-El-Haija, A. Kapoor, B. Perozzi, J. Lee, N-gcn: Multi-scale graph convolution for semi-supervised node classification, *arXiv preprint arXiv:1802.08888* (2018).
- [9] R. Al-Rfou, B. Perozzi, D. Zelle, Ddggk: Learning graph representations for deep divergence graph kernels, in: *Proceedings of the World Wide Web Conference*, 2019, pp. 37–48.
- [10] R. C. Fong, A. Vedaldi, Interpretable explanations of black boxes by meaningful perturbation, in: *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 3429–3437.

- [11] Y. Niu, L. Gu, F. Lu, F. Lv, Z. Wang, I. Sato, Z. Zhang, Y. Xiao, X. Dai, T. Cheng, Pathological evidence exploration in deep retinal image diagnosis, in: Proceedings of the AAAI conference on Artificial Intelligence, 2019, pp. 1093–1101.
- [12] H. Lakkaraju, E. Kamar, R. Caruana, J. Leskovec, Interpretable & explorable approximations of black box models, arXiv preprint arXiv:1707.01154 (2017).
- [13] J. Wagner, J. M. Kohler, T. Gindele, L. Hetzel, J. T. Wiedemer, S. Behnke, Interpretable and fine-grained visual explanations for convolutional neural networks, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2019, pp. 9097–9107.
- [14] M. T. Ribeiro, S. Singh, C. Guestrin, Why should i trust you?: Explaining the predictions of any classifier, in: Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2016, pp. 1135–1144.
- [15] Z. Ying, D. Bourgeois, J. You, M. Zitnik, J. Leskovec, Gnnexplainer: Generating explanations for graph neural networks, in: Proceedings of the Advances in Neural Information Processing Systems, 2019, pp. 9240–9251.
- [16] Z. Chen, X. Wei, P. Wang, Y. Guo, Multi-label image recognition with graph convolutional networks, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2019, pp. 5177–5186.
- [17] L. Shi, Y. Zhang, J. Cheng, H. Lu, Two-stream adaptive graph convolutional networks for skeleton-based action recognition, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2019, pp. 12026–12035.
- [18] L. Yao, C. Mao, Y. Luo, Graph convolutional networks for text classification, in: Proceedings of the AAAI Conference on Artificial Intelligence, 2019, pp. 7370–7377.
- [19] R. Ying, R. He, K. Chen, P. Eksombatchai, W. L. Hamilton, J. Leskovec, Graph convolutional neural networks for web-scale recommender systems, in: Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2018, pp. 974–983.

- [20] J. Shang, C. Xiao, T. Ma, H. Li, J. Sun, Gamenet: graph augmented memory networks for recommending medication combination, in: Proceedings of the AAAI Conference on Artificial Intelligence, 2019, pp. 1126–1133.
- [21] Z. Zhang, P. Cui, W. Zhu, Deep learning on graphs: A survey, *IEEE Trans. Knowl. Data Eng.* (2020).
- [22] Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, S. Y. Philip, A comprehensive survey on graph neural networks, *IEEE Trans. Neural Netw. Learn. Syst.* (2020).
- [23] R. Li, S. Wang, F. Zhu, J. Huang, Adaptive graph convolutional neural networks, in: Proceedings of the AAAI conference on artificial intelligence, 2018, pp. 3546–3553.
- [24] C. Zhuang, Q. Ma, Dual graph convolutional networks for graph-based semi-supervised classification, in: Proceedings of the World Wide Web Conference, 2018, pp. 499–508.
- [25] K. Xu, W. Hu, J. Leskovec, S. Jegelka, How powerful are graph neural networks?, *arXiv preprint arXiv:1810.00826* (2018).
- [26] M. Zhang, Z. Cui, M. Neumann, Y. Chen, An end-to-end deep learning architecture for graph classification, in: Proceedings of the AAAI Conference on Artificial Intelligence, 2018, pp. 4438–4445.
- [27] W. Chiang, X. Liu, S. Si, Y. Li, S. Bengio, C. Hsieh, Cluster-gcn: An efficient algorithm for training deep and large graph convolutional networks, in: Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2019, pp. 257–266.
- [28] J. Bruna, W. Zaremba, A. Szlam, Y. LeCun, Spectral networks and locally connected networks on graphs, *arXiv preprint arXiv:1312.6203* (2013).
- [29] R. Levie, F. Monti, X. Bresson, M. M. Bronstein, Caylennets: Graph convolutional neural networks with complex rational spectral filters, *IEEE Trans. Signal Process.* 67 (1) (2018) 97–109.

- [30] Z. Ying, J. You, C. Morris, X. Ren, W. Hamilton, J. Leskovec, Hierarchical graph representation learning with differentiable pooling, in: Proceedings of the Advances in Neural Information Processing Systems, 2018, pp. 4800–4810.
- [31] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Lio, Y. Bengio, Graph attention networks, arXiv preprint arXiv:1710.10903 (2017).
- [32] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, A. Torralba, Learning deep features for discriminative localization, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016, pp. 2921–2929.
- [33] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, D. Batra, Grad-cam: Visual explanations from deep networks via gradient-based localization, in: Proceedings of the IEEE International Conference on Computer Vision, 2017, pp. 618–626.
- [34] A. Chattopadhyay, A. Sarkar, P. Howlader, V. N. Balasubramanian, Grad-cam++: Generalized gradient-based visual explanations for deep convolutional networks, in: Proceedings of the IEEE Winter Conference on Applications of Computer Vision, 2018, pp. 839–847.
- [35] J. Zhang, S. A. Bargal, Z. Lin, J. Brandt, X. Shen, S. Sclaroff, Top-down neural attention by excitation backprop, *Int. J. Comput. Vis.* 126 (10) (2018) 1084–1102.
- [36] S. Bach, A. Binder, G. Montavon, F. Klauschen, K.-R. Müller, W. Samek, On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation, *PLoS One* 10 (7) (2015) e0130140.
- [37] A. Mahendran, A. Vedaldi, Salient deconvolutional networks, in: Proceedings of the European Conference on Computer Vision, 2016, pp. 120–135.
- [38] S. M. Lundberg, S.-I. Lee, A unified approach to interpreting model predictions, in: Proceedings of the Advances in Neural Information Processing Systems, 2017, pp. 4765–4774.
- [39] P. E. Pope, S. Kolouri, M. Rostami, C. E. Martin, H. Hoffmann, Explainability methods for graph convolutional neural networks, in: Proceedings

of the IEEE Conference on Computer Vision and Pattern Recognition, 2019, pp. 10772–10781.

- [40] F. Baldassarre, H. Azizpour, Explainability techniques for graph convolutional networks, arXiv preprint arXiv:1905.13686 (2019).
- [41] A. K. Debnath, R. L. Lopez de Compadre, G. Debnath, A. J. Shusterman, C. Hansch, Structure-activity relationship of mutagenic aromatic and heteroaromatic nitro compounds. correlation with molecular orbital energies and hydrophobicity, *J. Med. Chem.* 34 (2) (1991) 786–797.
- [42] P. Yanardag, S. Vishwanathan, Deep graph kernels, in: Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2015, pp. 1365–1374.