
GOAt: Explaining Graph Neural Networks via Graph Output Attribution

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Understanding the decision-making process of Graph Neural Networks (GNNs)
2 is crucial to their interpretability. Present methods for explaining GNNs typically
3 rely on training auxiliary models, and may struggle with issues such as overfitting
4 to noise, insufficient discriminability, and inconsistent explanations across data
5 samples of the same class. This paper introduces Graph Output Attribution (GOAt),
6 a novel method to attribute graph outputs to input graph features, creating GNN
7 explanations that are faithful, discriminative, as well as stable across similar sam-
8 ples. By expanding the GNN as a sum of scalar products involving node features,
9 edge features and activation patterns, we propose an efficient analytical method
10 to compute contribution of each node or edge feature to each scalar product and
11 aggregate the contributions from all scalar products in the expansion form to derive
12 the importance of each node and edge. Through extensive experiments on synthetic
13 and real data, we show that our method has consistently outperformed various
14 state-of-the-art GNN explainers in terms of fidelity, discriminability, and stability.

15 1 Introduction

16 Graph Neural Networks (GNNs) have demonstrated notable success in learning representations from
17 graph-structured data in various fields [14, 9, 30]. However, their black-box nature has driven the
18 need for explainability, especially in sectors where transparency and accountability are essential, such
19 as finance [28], healthcare [1], and security [18]. The ability to interpret GNNs can provide insights
20 into the mechanisms underlying deep models and help establish trustworthy predictions.

21 Existing attempts to explain GNNs usually focus on local-level or global-level explainability. Local-
22 level explainers [31, 17, 21, 27, 10, 15, 23, 3] typically train a secondary model to identify the critical
23 graph structures that best explain the behavior of a pretrained GNN for specific input instances. These
24 methods are always optimized for ground-truth explanations or fidelity metrics, yet may not be able to
25 generate consistent explanations for similar graph samples or produce accurate and human-intelligible
26 explanations for class discrimination. Global-level explainers [2, 11] perform prototype learning or
27 random walk on the explanation instances to extract the global explanations over a multitude of graph
28 samples. However, their effectiveness rely heavily on the quality of local-level explanations.

29 In this paper, we introduce a computationally efficient local-level GNN explanation technique called
30 **Graph Output Attribution (GOAt)** to overcome the limitations of existing methods. Unlike methods
31 that rely on back-propagation with gradients [20, 4, 22, 8] and those relying on hyper-parameters or
32 training complex black-box models [17, 15, 23, 3], our approach enables attribution of GNN output
33 to input features, leveraging the repetitive sum-product structure in the forward pass of a GNN.

34 Given that the matrix multiplication in each GNN layer adheres to linearity properties and the
35 activation functions operate element-wise, a GNN can be represented in an expansion form as a

36 sum of scalar product terms, involving input graph features, model parameters, as well as *activation*
 37 *patterns* that indicate the activation levels of the scalar products. Based on the notion that all scalar
 38 variables X_i in a scalar product term $g = cX_1X_2 \dots X_N$ contribute equally to g , where c is a
 39 constant, we can attribute each product term to its corresponding factors and thus to input features,
 40 obtaining the importance of each node or edge feature in the input graph to GNN outputs. We present
 41 case studies that demonstrate the effectiveness of our analytical explanation method *GOAt* on typical
 42 GNN variants, including GCN, GraphSAGE, and GIN.

43 Besides the fidelity metric, which is commonly used to assess the faithfulness of GNN explanations,
 44 we introduce two new metrics to evaluate the *discriminability* and *stability* of the explanation, which
 45 are under-investigated by prior literature. Discriminability refers to the ability of explanations to
 46 distinguish between classes, which is assessed by the difference between the mean explanation
 47 embeddings of different classes, while stability refers to the ability to generate consistent explanations
 48 across similar data instances, which is measured by the percentage of data samples covered by top- k
 49 explanations. Through comprehensive experiments based on both synthetic and real-world datasets
 50 along with qualitative analysis, we show the outstanding performance of our proposed method, *GOAt*,
 51 in providing highly faithful, discriminative, and stable explanations for GNNs, as compared to a
 52 range of state-of-the-art methods.

53 2 Problem Formulation

54 **Graph Neural Networks** Let $G = (\mathcal{V}, \mathcal{E})$ be a graph, where $\mathcal{V} = \{v_1, v_2, \dots, v_N\}$ denotes the set
 55 of nodes and $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ denotes the set of edges. The node feature matrix of the graph is represented
 56 by $X \in \mathbb{R}^{N \times d}$, and the adjacency matrix is represented by $A \in \{0, 1\}^{N \times N}$ such that $A_{ij} = 1$ if
 57 there exists an edge between nodes v_i and v_j . The task of a GNN is to learn a function $f(G)$, which
 58 maps the input graph G to a target output, such as node labels, graph labels, or edge labels. Formally
 59 speaking, for a given GNN, the hidden state $h_i^{(l)}$ of node v_i at its layer l can be represented as:

$$h_i^{(l)} = \text{COMBINE}^{(l)} \left\{ h_i^{(l-1)}, \text{AGGREGATE}^{(l)} \left(\left\{ h_j^{(l-1)}, \forall v_j \in \mathcal{N}_i \right\} \right) \right\}, \quad (1)$$

60 where \mathcal{N}_i represents the set of neighbors of node v_i in the graph. COMEBINE^(l)(\cdot) is a COMBINE
 61 function such as concatenation [9], while AGGREGATE^(l)(\cdot) are AGGREGATE functions with
 62 aggregators such as ADD. We focus on GNNs that adopt the non-linear activation function ReLU in
 63 COMBINE or AGGREGATE functions.

64 **Local-level GNN Explainability** Our goal is to generate a faithful explanation for a graph instance
 65 $G = (\mathcal{V}, \mathcal{E})$ by identifying a subset of edges $S \subseteq \mathcal{E}$, given a GNN $f(\cdot)$ pretrained on a set of graphs
 66 \mathcal{G} . The term *faithful* refers to the explanation’s ability to perform well in not only *fidelity* [33] and
 67 *robustness* [3] metrics, but also *stability* in identifying consistent patterns. We highlight edges instead
 68 of nodes as suggested by [7] that edges have more fine-grained information than nodes while giving
 69 human-understandable explanations like subgraphs.

70 3 Method

71 This section begins by presenting our fundamental definition of equal contribution in a product term
 72 and its application in an example of a toy graph neural network. Then, we mathematically present
 73 *GOAt* method for explaining typical GNNs, followed by a case study on GCN [14]. Additional case
 74 studies of applying *GOAt* to GraphSAGE [9] and GIN [30] are included in the Appendix.

75 3.1 Definitions

76 Consider a function $g(X_1, \dots, X_M)$ of M variables $\mathbf{X} = \{X_1, \dots, X_M\}$. If we let a pair of variables
 77 (X_i, X_j) be set to $(X_i, X_j) = (x_i, x_j)$, we will obtain a manifold $g_{X_i=x_i, X_j=x_j}(\mathbf{X} \setminus \{X_i, X_j\})$,
 78 which represents $g(\cdot)$ when all variables excluding X_i and X_j can vary. Consider a base manifold
 79 $g_{X_i=x'_i, X_j=x'_j}(\mathbf{X} \setminus \{X_i, X_j\})$. If we can obtain two identical manifolds by setting $(X_i, X_j) =$
 80 (x_i, x'_j) and $(X_i, X_j) = (x'_i, x_j)$, it will indicate that changing $X_i = x'_i$ to $X_i = x_i$ is equivalent
 81 to changing $X_j = x'_j$ to $X_j = x_j$ with respect to the base manifold at $(X_i, X_j) = (x'_i, x'_j)$.
 82 For example, a function $g(x, y, z) = 2xy + x^2z$ has three variables x, y, z , we consider taking

83 $g_{x=-1,y=-1}(z) = z + 2$ as the base manifold. Since $g_{x=1,y=-1}(z) = g_{x=-1,y=1}(z) = z - 2$, we
 84 say that changing $x = -1$ to $x = 1$ is equivalent to changing $y = -1$ to $y = 1$ with respect to the
 85 base manifold at $(x, y) = (-1, -1)$.

86 **Definition 1 (Equal Contribution).** Given a function $g(\mathbf{X})$ where $\mathbf{X} = \{X_1, \dots, X_M\}$ represents
 87 M variables, we say that variables X_i and X_j have equal contribution to function g at (x_i, x_j)
 88 with respect to the base manifold at (x'_i, x'_j) if and only if setting $X_i = x_i, X_j = x'_j$ and setting
 89 $X_i = x'_i, X_j = x_j$ yield the same manifold, i.e.,

$$g_{X_i=x_i, X_j=x'_j}(\mathbf{X} \setminus \{X_i, X_j\}) = g_{X_i=x'_i, X_j=x_j}(\mathbf{X} \setminus \{X_i, X_j\})$$

90 for any values of \mathbf{X} excluding X_i and X_j .

91 **Lemma 2 (Equal Contribution in a product).** Given a function $g(\mathbf{X})$ defined as $g(\mathbf{X}) =$
 92 $b \prod_{k=1}^M X_k$, where b is a constant, and $\mathbf{X} = \{X_1, \dots, X_M\}$ represents M uncorrelated variables.
 93 Each variable X_k is either 0 or x_k , depending on the absence or presence of a certain feature. Then,
 94 all the variables in \mathbf{X} contribute equally to $g(\mathbf{X})$ at $[x_1, \dots, x_M]$ with respect to $[0, \dots, 0]$.

95 Proofs of all Lemmas and Theorems can be found in the Appendix. Since all the binary variables
 96 have equal contribution, we define the contribution of each variable X_k to $g(\mathbf{X}) = b \prod_{k=1}^M X_k$ for
 97 all $k = 1, \dots, M$, as

$$I_{X_k} = \frac{b \prod_{i=1}^M x_i}{M}. \quad (2)$$

98 For example, let $f(A, X) = AXW$ be a simple 2-node GNN for node classification, where A, X, W
 99 are 2×2 matrices that denote adjacency matrix, node feature matrix, and weight matrix, respectively.
 100 Then, we can represent each entry in the resulting 2×2 matrix $f(A, X)$ as an expansion form:

$$f_{i,j}(A, X) = \sum_{k=0}^1 \sum_{l=0}^1 A_{i,k} X_{k,l} W_{l,j}, \quad (3)$$

101 where $f_{i,j}(A, X)$ represents the prediction of the i -th node for the j -th class. In a pretrained
 102 GNN, parameter W is fixed. Thus, only $A_{i,k}$ and $X_{k,l}$ contribute to the value of each *scalar product*
 103 $A_{i,k} X_{k,l} W_{l,j}$. As $A_{i,k}$ is usually independent of $X_{k,l}$ under proper data cleaning, we can calculate the
 104 contributions of $A_{i,k}$ and $X_{k,l}$ to the scalar product $A_{i,k} X_{k,l} W_{l,j}$ by $I_{A_{i,k}} = I_{X_{k,l}} = \frac{1}{2} A_{i,k} X_{k,l} W_{l,j}$
 105 based on Lemma 2 and Equation (2). By similar computations for all the scalar products in the
 106 expansion form of $f(\cdot)$, we can obtain the contribution of all the input features to each entry of the
 107 output matrix.

108 3.2 Explaining Graph Neural Networks via Attribution

109 A typical GNN [14, 9, 30] for node or graph classification tasks usually comprises 2-6 message-
 110 passing layers for learning node or graph representations, followed by several fully connected layers
 111 that serve as the classifier. With the hidden state $h_i^{(l)}$ of node v_i at the l -th message-passing layer
 112 defined as Equation (1), we generally have the hidden state $H^{(l)}$ of a data sample as:

$$H^{(l)} = \sigma \left(\Phi^{(l)} \left((A + \epsilon^{(l)} I) H^{(l-1)} \right) + \lambda \Psi^{(l)} \left(H^{(l-1)} \right) \right), \quad (4)$$

113 where A is the adjacency matrix, $\epsilon^{(l)}$ refers to the self-loop added to the graph if fixed to 1, otherwise it
 114 is a learnable parameter, $\sigma(\cdot)$ is the element wise activation function, $\Phi^{(l)}$ and $\Psi^{(l)}$ can be Multilayer
 115 Perceptrons (MLP) or linear mappings, $\lambda \in \{0, 1\}$ determines whether a concatenation is required.
 116 If the COMBINE step of a GNN requires a concatenation, we have $\lambda = 1$ and $\epsilon^{(l)} = 1$; if the
 117 COMBINE step requires a weighted sum, we have $\epsilon^{(l)}$ set trainable and $\lambda = 0$. Alternatively,
 118 Equation (4) can be expanded to:

$$H^{(l)} = \sigma \left(AH^{(l-1)} \prod_{k=1}^K W^{\Phi_k^{(l)}} + \epsilon^{(l)} H^{(l-1)} \prod_{k=1}^K W^{\Phi_k^{(l)}} + \lambda H^{(l-1)} \prod_{q=1}^Q W^{\Psi_q^{(l)}} \right), \quad (5)$$

119 where K, Q refer to the number of MLP layers in $\Phi^{(l)}(\cdot)$ and $\Psi^{(l)}(\cdot)$, and $W^{\Phi_k^{(l)}}$ and $W^{\Psi_q^{(l)}}$ are the
 120 trainable parameters in $\Phi_k^{(l)}$ and $\Psi_q^{(l)}$.

121 Given a certain data sample and a pretrained GNN, for an element-wise activation function we can
 122 define the activation pattern as the ratio between the output and input of the activation function:

123 **Definition 3 (Activation Pattern).** Denote $H^{(l)'}$ and $H^{(l)}$ as the hidden representations before and
 124 after passing through the element-wise activation function at the l -th layer, we define activation
 125 pattern $P^{(l)}$ for a given data sample as

$$P_{i,j}^{(l)} = \begin{cases} \frac{H_{i,j}^{(l)}}{H_{i,j}^{(l)'}}, & \text{if } H_{i,j}^{(l)' } \neq 0 \\ 0, & \text{otherwise} \end{cases}$$

126 where $P_{i,j}^{(l)}$ is the element-wise activation pattern for the j -th feature of i -th node at layer l .

127 Hence, the hidden state $H^{(l)}$ at the l -th layer for a given sample can be written as

$$H^{(l)} = P^{(l)} \odot \left(AH^{(l-1)} \prod_{k=1}^K W^{\Phi_k^{(l)}} + \epsilon^{(l)} H^{(l-1)} \prod_{k=1}^K W^{\Phi_k^{(l)}} + \lambda H^{(l-1)} \prod_{q=1}^Q W^{\Psi_q^{(l)}} \right), \quad (6)$$

128 where \odot represents element-wise multiplication. Thus, similar to Equation (3), we can expand the
 129 expression of each output entry in a GNN $f(A, X)$ into a sum of scalar products, where each scalar
 130 product is the multiplication of corresponding entries in A , X , W , and P in all layers. Then each
 131 scalar product can be written as

$$z = \mathbb{C} \cdot \left(P_{\alpha_{10}, \beta_{11}}^{(1)} \cdots P_{\alpha_{L0}, \beta_{L1}}^{(L)} \right) \left(P_{\alpha_{L0}, \gamma_{11}}^{(c_1)} \cdots P_{\alpha_{L0}, \gamma_{(M-1)1}}^{(c_{M-1})} \right) \cdot \\ \left(A_{\alpha_{L0}, \alpha_{L1}}^{(L)} \cdots A_{\alpha_{10}, \alpha_{11}}^{(1)} \right) X_{i,j} \left(W_{\beta_{10}, \beta_{11}}^{(1)} \cdots W_{\beta_{L0}, \beta_{L1}}^{(L)} \right) \left(W_{\gamma_{10}, \gamma_{11}}^{(c_1)} \cdots W_{\gamma_{M0}, \gamma_{M1}}^{(c_M)} \right), \quad (7)$$

132 where \mathbb{C} is a constant, c_k refers to the k -th layer of the classifier, $(\alpha_{l0}, \alpha_{l1})$, (β_{l0}, β_{l1}) , $(\gamma_{l0}, \gamma_{l1})$ are
 133 (row, column) indices of the corresponding matrices at layer l . In a classifier with M MLP layers,
 134 only $(M - 1)$ layers adopt activation functions. Therefore, we do not have $P_{\alpha_{L0}, \gamma_{M1}}^{(c_M)}$ in Equation (7).
 135 For scalar products without factors of A , all A 's are considered as constants equal to 1 in Equation (7).
 136 Since the GNN model parameters are pretrained and fixed, we only consider A , X , and all the P
 137 terms as the variables in each product term.

138 **Lemma 4 (Equal Contribution variables in the GNN expansion form's scalar product).** For a
 139 scalar product term z in the expansion form of a pretrained GNN $f(\cdot)$, when the number of nodes N
 140 is large, all variables in z have equal contributions to the scalar product z .

141 Hence, by Equation (2), we can calculate the contribution $I_\nu(z)$ of a variable ν (i.e., an entry in A , X
 142 and P matrices) to each scalar product z (given by Equation (7)) by:

$$I_\nu(z) = \frac{z}{|V(z)|}, \quad (8)$$

143 where function $V(\cdot)$ represents the set of variables in its input, and $|V(z)|$ denotes the number of
 144 unique **variables** in z , e.g., $V(x^2y) = \{x, y\}$, and $|V(x^2y)| = 2$.

145 Similar to Section 3.1, an entry $f_{m,n}(A, X)$ of the output matrix $f(A, X)$ can be expressed by the
 146 sum of all the related scalar products as

$$f_{m,n}(A, X) = \sum \mathbb{C} \cdot \left(P_{\alpha_{10}, \beta_{11}}^{(1)} \cdots P_{\alpha_{m}, \beta_{L1}}^{(L)} \right) \cdot \left(P_{\alpha_{m}, \gamma_{11}}^{(c_1)} \cdots P_{\alpha_{m}, \gamma_{(M-1)1}}^{(c_{M-1})} \right) \cdot \left(A_{\alpha_{m}, \alpha_{L1}}^{(L)} \cdots A_{\alpha_{10}, \alpha_{11}}^{(1)} \right) \\ \cdot X_{i,j} \cdot \left(W_{\beta_{10}, \beta_{11}}^{(1)} \cdots W_{\beta_{L0}, \beta_{L1}}^{(L)} \right) \cdot \left(W_{\gamma_{10}, \gamma_{11}}^{(c_1)} \cdots W_{\gamma_{M0}, \gamma_{M1}}^{(c_M)} \right), \quad (9)$$

147 where summation is over all possible $(\alpha_{l0}, \alpha_{l1})$, (β_{l0}, β_{l1}) , $(\gamma_{l0}, \gamma_{l1})$, for message passing layer
 148 $l = 1, \dots, L$ or classifier layer $l = 1, \dots, M$, as well as all i, j indices for X . By summing up the
 149 contribution of each variable ν among the entries in the A , X and P 's in all the scalar products in the
 150 expansion form of $f_{m,n}(\cdot)$, we can obtain the contribution of ν to $f_{m,n}(\cdot)$ as:

$$I_\nu(f_{m,n}(\cdot)) = \sum_{z \text{ in } f_{m,n}(\cdot) \text{ that contain } \nu} \frac{z}{|V(z)|}. \quad (10)$$

151 **Theorem 5 (Contribution of variables in the expansion form of a pretrained GNN).** Given
 152 Equations (8) and (10), for each variable ν (i.e., an entry in A , X and P matrices), when the number
 153 of nodes N is large, we can approximate $I_\nu(f_{m,n}(\cdot))$ by:

$$I_\nu(f_{m,n}(\cdot)) = \sum_{z \text{ in } f_{m,n}(\cdot) \text{ that contain } \nu} \frac{O(\nu, z)}{\sum_{\rho \text{ in } z} O(\rho, z)} \cdot z, \quad (11)$$

154 where $O(\nu, z)$ denotes the number of occurrences of ν among the variables of z .

155 Recall that $|V(z)|$ stand for the number of **unique variables** in z . Hence the total number of
 156 occurrences of all the variables $\sum_{\rho \text{ in } z} O(\rho, z)$ is not necessarily equal to $|V(z)|$. For example,
 157 if all of $\{A_{\alpha_{10}, \alpha_{11}}^{(1)}, \dots, A_{\alpha_{L0}, \alpha_{L1}}^{(L)}\}$ in z are unique entries in A , then they can be considered as L
 158 independent variables in the function representing z . If two of these occurrences of variables refer to
 159 the same entry in A , then there are only $(L - 1)$ unique variables related to A .

160 Although Theorem 5 gives the contribution of each entry in A , X and P 's, we need to further attribute
 161 P 's to A and X and allocate the contribution of each activation pattern $P_{a,b}^{(r)}$ to node features X and
 162 edges A by considering all non-zero features in X_a of node v_a and the edges within m hops of node
 163 v_a , as these inputs may contribute to the activation pattern $P_{a,b}^{(r)}$. However, determining the exact
 164 contribution of each feature that contributes to $P_{a,b}^{(r)}$ is not straightforward due to non-linear activation.
 165 We approximately attribute all relevant features equally to $P_{a,b}^{(r)}$. That is, each input feature ν that has
 166 nonzero contribution to $P_{a,b}^{(r)}$ will share an equal contribution of $I_{P_{a,b}^{(r)}}(f_{m,n}(\cdot))/|V(P_{a,b}^{(r)})|$, where
 167 $|V(P_{a,b}^{(r)})|$ denotes the number of distinct node and edge features in X and A contributing to $P_{a,b}^{(r)}$,
 168 which is exactly all non-zero features in X_a of node v_a and the adjacency matrix entries within r
 169 hops of node v_a . Finally, based on Equation (11), we can obtain the contribution of an input feature ν
 170 in X , A of a graph instance to the (m, n) -th entry of the GNN output $f(\cdot)$ as:

$$\hat{I}_\nu(f_{m,n}(\cdot)) = I_\nu(f_{m,n}(\cdot)) + \sum_{P_{a,b}^{(r)} \text{ in } f_{m,n}(\cdot), \text{ with } \nu \text{ in } P_{a,b}^{(r)}} \frac{I_{P_{a,b}^{(r)}}(f_{m,n}(\cdot))}{|V(P_{a,b}^{(r)})|}, \quad (12)$$

171 where ν is an entry in the adjacency matrix A or the input feature matrix X , $P_{a,b}^{(r)}$ denotes an entry in
 172 all the activation patterns. Thus, we have attributed $f(\cdot)$ to each input feature of a given data instance.

173 Our approach meets the *completeness* axiom, which is a critical requirement in attribution methods [25,
 174 24, 6]. This axiom guarantees that the attribution scores for input features add up to the difference
 175 in the GNN's output with and without those features. Passing this sanity check implies that our
 176 approach provides a more comprehensive account of feature importance than existing methods that
 177 only rank the top features [3, 17, 20, 31, 27, 23].

178 3.3 Case Study: Explaining Graph Convolutional Network (GCN)

179 GCNs use a simple sum in the combination step, and the adjacency matrix is normalized with the
 180 diagonal node degree matrix D . Hence, the hidden state of a GCN's l -th message-passing layer is:

$$H^{(l)} = \text{ReLU} \left(V H^{(l-1)} W^{(l)} + B^{(l)} \right), \quad (13)$$

181 where $V = D^{-\frac{1}{2}}(A + I)D^{-\frac{1}{2}}$ represents the normalized adjacency matrix with self-loops added.
 182 Suppose a GCN has three convolution layers and a 2-layer MLP as the classifier, then its expansion
 183 form without the activation functions $\text{ReLU}(\cdot)$ will be:

$$\begin{aligned} f(V, X)_{\mathcal{P}} &= V^{(3)} V^{(2)} V^{(1)} X W^{(1)} W^{(2)} W^{(3)} W^{(c_1)} W^{(c_2)} + V^{(3)} V^{(2)} B^{(1)} W^{(2)} W^{(3)} W^{(c_1)} W^{(c_2)} \\ &+ V^{(3)} B^{(2)} W^{(3)} W^{(c_1)} W^{(c_2)} + B^{(3)} W^{(c_1)} W^{(c_2)} + B^{(c_1)} W^{(c_2)} + B^{(c_2)}, \end{aligned} \quad (14)$$

184 where $V^{(l)} = V$ is the normalized adjacency matrix in the l -th layer's calculation. In the actual
 185 expansion form with the activation patterns, the corresponding $P^{(m)}$'s are multiplied whenever
 186 there is a $W^{(m)}$ or $B^{(m)}$ in a scalar product, excluding the last layer $W^{(c_2)}$ and $B^{(c_2)}$. For ex-
 187 ample, in the scalar products corresponding to $V^{(3)} V^{(2)} V^{(1)} X W^{(1)} W^{(2)} W^{(3)} W^{(c_1)} W^{(c_2)}$, there
 188 are eight variables consisting of four P 's, one X , and three V 's. By Equation (11), an ad-
 189 jacency entry $V_{i,j}$ itself will contribute $\frac{1}{8}$ of $p(V^{(3)} V_{i,j}^{(2)} V_{i,j}^{(1)} X_j : W^{(1)} W^{(2)} W^{(3)} W^{(c_1)} W^{(c_2)}) +$
 190 $p(V_{i,j}^{(3)} V_{i,j}^{(2)} V_j^{(1)} X W^{(1)} W^{(2)} W^{(3)} W^{(c_1)} W^{(c_2)}) + p(V_{i,j}^{(3)} V_j^{(2)} V^{(1)} X W^{(1)} W^{(2)} W^{(3)} W^{(c_1)} W^{(c_2)})$,
 191 where $p(\cdot)$ denotes the results with the element-wise multiplication of the corresponding activation
 192 patterns applied at the appropriate layers. After we obtain the contribution of $V_{i,j}$ itself on all the
 193 scalar products, we can follow Equation (12) to allocate the contribution of activation patterns to $V_{i,j}$.

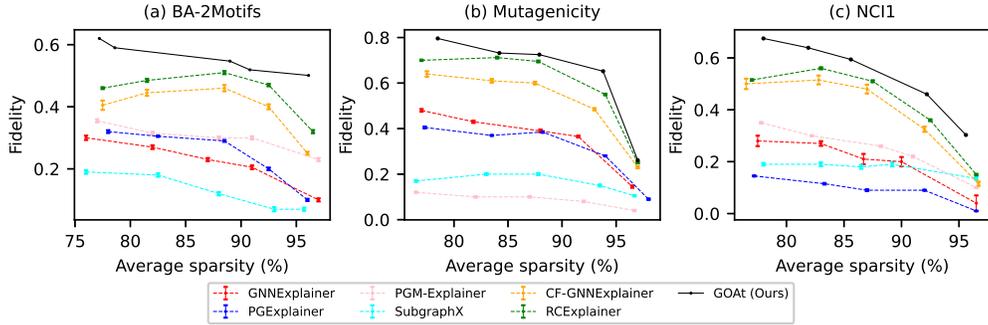


Figure 1: Fidelity performance averaged across 10 runs on the pretrained GCNs for the datasets at different levels of average sparsity.

194 With Equation (14), we find that when both V and X are set to zeros, $f(\cdot)$ remains non-zero and is:

$$f(\mathbf{0}, \mathbf{0}) = p(B^{(3)}W^{(c_1)}W^{(c_2)}) + p(B^{(c_1)}W^{(c_2)}) + B^{(c_2)}, \quad (15)$$

195 where $B^{(c_2)}$ is the global bias, and the other terms have non-zero entries at the activated neurons. In
 196 other words, certain GNN neurons in the 3-rd and c_1 -th layers may already be activated prior to any
 197 input feature being passed to the network. When we do feed input features, some of these neurons
 198 may remain activated or be toggled off. With Equation (12), we consider taking all 0's of the X
 199 entries, V entries and P entries as the base manifold. Now, given that some of the P entries in GCN
 200 are non-zero when all X and V set to zeros, as present in Equation (15), we will need to subtract the
 201 contribution of each features on these P from the contribution values calculated by Equation (12).
 202 We let \mathbf{P}' represent the activation patterns of $f(\mathbf{0}, \mathbf{0})$, then the calibrated contribution $\hat{I}_{V_{i,j}}^{\text{cali}}(f(\cdot))$ of
 203 $V_{i,j}$ is given by:

$$\hat{I}_{V_{i,j}}^{\text{cali}}(f(\cdot)) = \hat{I}_{V_{i,j}}(f(V, X)) - \sum_{P'_{a,b} \text{ in } f(\mathbf{0}, \mathbf{0}), \text{ with } V_{i,j} \text{ in } P'_{a,b}} \frac{I_{P'_{a,b}}(f(\mathbf{0}, \mathbf{0}))}{|V(P'_{a,b})|}. \quad (16)$$

204 In graph classification tasks, a pooling layer such as mean-pooling is added after the convolution
 205 layers to obtain the graph representation. To determine the contribution of each input feature, we can
 206 simply apply the same pooling operation as used in the pre-trained GCN.

207 As we mentioned in Section 2, we aim to obtain the explanations by the critical edges in this paper,
 208 since edges have more fine-grained information than nodes. Therefore, we treat the edges as variables,
 209 while considering the node features X as constants similar to parameters W or B . This setup naturally
 210 aggregates the contribution of node features onto edges. By leveraging edge attributions, we are able
 211 to effectively highlight motifs within the graph structure.

212 4 Experiments

213 We conduct a series of experiments on the fidelity, discriminability and stability metrics to compare
 214 our method with the state-of-the-art methods including GNNExplainer [31], PGExplainer [17], PGM-
 215 Explainer [27], SubgraphX [33], CF-GNNExplainer [16], RCEExplainer [3], RG-Explainer [23] and
 216 DEGREE [8]. As outlined in Section 2, we highlight edges as explanations as suggested by [7]. For
 217 baselines that identify nodes or subgraphs as explanations, we adopt the evaluation setup from [3].

218 We evaluate the performance of explanations on three variants of GNNs, which are GCN [14],
 219 GraphSAGE [9] and GIN [30]. The experiments are conducted on both the graph classification
 220 task and the node classification task. For graph classification task, we evaluate on a synthetic
 221 dataset, BA-2motifs [17], and two real-world datasets, Mutagenicity [13] and NCI1 [19]. For node
 222 classification task, we evaluate on three synthetic datasets [17], which are BA-shapes, BA-Community
 223 and Tree-grid. As space is limited, we will only present the key results here. Fidelity results on GIN
 224 and GraphSAGE, as well as the results of node classification tasks can be found in the Appendix.
 225 Discussions on the controversial metrics such as accuracy are also moved to the Appendix.

226 4.1 Fidelity

227 Fidelity [20, 32, 29, 3] is the decrease of predicted probability between original and new predictions
 228 after removing important edges, which are used to evaluate the faithfulness of explanations. It is

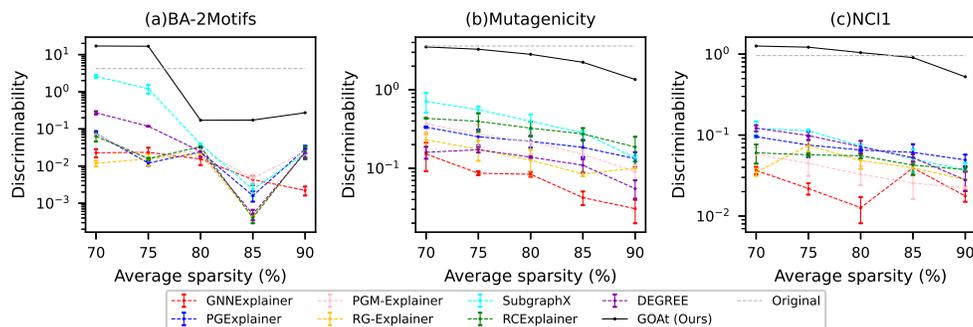


Figure 2: Discriminability performance averaged across 10 runs of the explanations produced by various GNN explainers at different levels of sparsity. "Original" refer to the performance of feeding the original data into the GNN without any modifications or explanations applied.

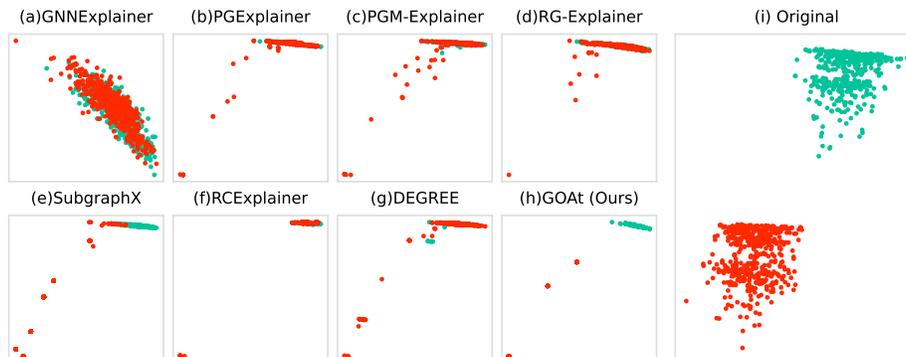


Figure 3: Visualization of explanation embeddings on the BA-2Motifs dataset. Subfigure (i) refers to the visualization of the original embeddings by directly feeding the original data into the GNN without any modifications or explanations applied.

229 defined as $fidelity(S, G) = f_y(G) - f_y(G \setminus S)$. As pointed out by [32], the fidelity may be sensitive
 230 to sparsity of explanations. The sparsity of an explanation $S \subseteq \mathcal{E}$ for a graph $G = \{\mathcal{V}, \mathcal{E}\}$ is given by
 231 $sparsity(S, G) = 1 - \frac{|S|}{|\mathcal{E}|}$. It indicates the percentage of edges that remain in G after the removal of
 232 edges in S . Higher sparsity means fewer edges are identified as critical, which may have a smaller
 233 impact on the prediction probability. Therefore, we compare fidelity performance under similar levels
 234 of average sparsity, as in [33, 29, 3]. Figure 1 displays the fidelity results, with the baseline results
 235 sourced from [3]. Our proposed approach, *GOAt*, consistently outperforms the baselines in terms
 236 of fidelity across all sparsity levels, validating its superior performance in generating accurate and
 237 reliable faithful explanations. Among the other methods, *RCExplainer* exhibits the highest fidelity, as
 238 it is specifically designed for fidelity optimization. Notably, unlike the other methods that require
 239 training and hyperparameter tuning, *GOAt* offers the advantage of being a training-free approach,
 240 thereby avoiding any errors across different runs.

241 4.2 Discriminability

242 Discriminability, also known as discrimination ability [5, 12], refers to the ability of the explanations
 243 to distinguish between the classes. We define the discriminability between two classes c_1 and c_2 as the
 244 L2 norm of the difference between the mean values of explanation embeddings of the two classes. The
 245 embeddings used for explanations are taken prior to the last-layer classifier, with node embeddings
 246 employed for node classification tasks and graph embeddings utilized for graph classification tasks.
 247 In this procedure, only the explanation subgraph S is fed into the GNN instead of G .

248 We show the discriminability across various sparsity levels on GCN, as illustrated in Figure 2. Due to
 249 the significant performance gap between the baselines and *GOAt*, a logarithmic scale is employed.
 250 Our approach consistently outperforms the baselines in terms of discriminability across all sparsity
 251 levels, demonstrating its superior ability to generate accurate and reliable class-specific explanations.
 252 Notably, at $sparsity = 0.7$, *GOAt* achieves higher discriminability than the original graphs on the
 253 BA-2Motifs and NCI1 datasets. This indicates that *GOAt* effectively reduces noise unrelated to
 254 the investigated class while producing informative class explanations. Additionally, we observe a

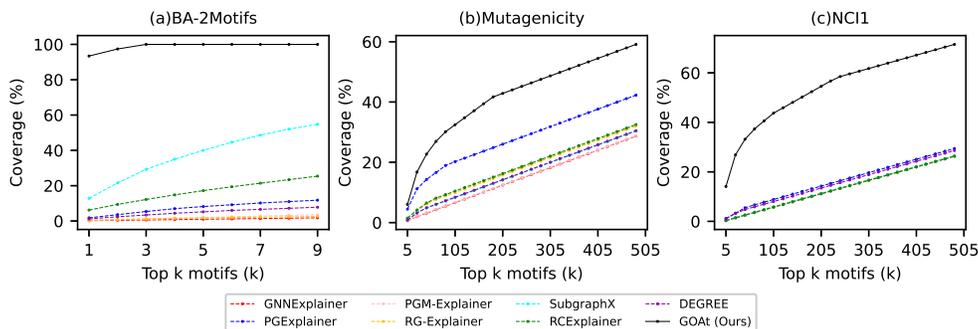


Figure 4: Coverage of the top- k explanations across the datasets.

255 substantial decrease in discriminability between sparsity levels of 0.75 and 0.8 on BA-2Motifs. This
 256 implies that a minimum of approximately 25% of the edges is necessary to distinguish between
 257 the classes, which is in line with our expectation, given that a "house" motif, consisting of 6 edges,
 258 usually represents 24% of the total edges (on average, the total number of edges in BA-2Motifs is 25).

259 Furthermore, we present scatter plots to visualize the explanation embeddings generated by various
 260 GNN explainers. Figure 3 showcases the explanation embeddings obtained from different GNN
 261 explaining methods on the BA-2Motifs dataset, with *sparsity* = 0.7. More scatter plots on Muta-
 262 genicity and NCI1 and can be found in the Appendix. The explanations generated by GNNExplainer
 263 fail to exhibit class discrimination, as all the data points are clustered together without any distinct
 264 separation. While some of the Class 1 explanations produced by PGExplainer, PGM-Explainer,
 265 RG-Explainer, RCEExplainer, and DEGREE are noticeably separate from the Class 0 explanations, the
 266 majority of the data points remain closely clustered together. As for SubgraphX, most of its Class 1
 267 explanations are isolated from the Class 0 explanations, but there is a discernible overlap between the
 268 Class 1 and Class 0 data points. In contrast, our method, *GOAt*, generates explanations that clearly
 269 and effectively distinguish between Class 0 and Class 1, with no overlapping points and a substantial
 270 separation distance, highlighting the strong discriminability of our approach.

271 4.3 Stability of extracting motifs

272 As we will later show in Section 4.4, it is often observed that datasets contain specific class motifs.
 273 For instance, in the BA-2Motifs dataset, the Class 1 motif exhibits a "house" structure. To ensure
 274 the stability of GNN explainers in capturing the class motifs across diverse data samples, we aim
 275 for the explanation motifs to exhibit relative consistency for data samples with similar properties,
 276 rather than exhibiting significant variations. To quantify this characteristic, we introduce the *stability*
 277 metric, which measures the coverage of the top- k explanations across the dataset. An ideal explainer
 278 should generate explanations that cover a larger number of data samples using fewer motifs. This
 279 characteristic is also highly desirable in global-level explainers, such as [2, 11]. We illustrate the
 280 stability of the unbiased class as the percentage converge of the top- k explanations produced on
 281 GCN with *sparsity* = 0.7 in Figure 4. Our approach surpasses the baselines by a considerable
 282 margin in terms of the stability of producing explanations. Specifically, *GOAt* is capable of providing
 283 explanations for all the Class 1 data samples using only three explanations. This explains why there
 284 are only three Class 1 scatters visible in Figure 3.

285 4.4 Qualitative analysis

286 We present the qualitative results of our approach in Table 1, where we compare it with state-of-the-art
 287 baselines such as PGExplainer, SubgraphX, and RCEExplainer. The pretrained GNN achieves a 100%
 288 accuracy on the BA-2Motifs dataset. As long as it successfully identifies one class, the remaining
 289 data samples naturally belong to the other class, leading to a perfect accuracy rate. Based on the
 290 explanations from *GOAt*, we have observed that the GNN effectively recognizes the "house" motif
 291 that is associated with Class 1. In contrast, other approaches face difficulties in consistently capturing
 292 this motif. The Class 0 motifs in the Mutagenicity dataset generated by *GOAt* represent multiple
 293 connected carbon rings. This indicates that the presence of more carbon rings in a molecule increases
 294 its likelihood of being mutagenic (Class 0), while the presence of more "C-H" or "O-H" bonds in a
 295 molecule increases its likelihood of being non-mutagenic (Class 1). Similarly, in the NCI1 dataset,
 296 *GOAt* discovers that the GNN considers a higher number of carbon rings as evidence of chemical

Table 1: Qualitative results of the top motifs of each class produced by PGExplainer, SubgraphX, RCEExplainer and *GOAt*. The percentages indicate the coverage of the explanations.

	BA-2Motifs		Mutagenicity		NCII							
	Class0	Class1	Class0	Class1	Class0	Class1						
PGExplainer	 4.8%	 1.8%	 1.2%	 1.3%	 0.1%	 0.5%						
SubgraphX	 0.4%	 12.8%	 0.2%	 0.2%	 0.2%	 0.1%						
RCEExplainer	 6.4%	 6.2%	 0.4%	 0.5%	 0.05%	 0.1%						
<i>GOAt</i>	 3.8%	 3.4%	 93.4%	 4%	 3.5%	 2.2%	 2.2%	 1.2%	 3.5%	 1.2%	 4.3%	 4.0%

297 compounds being active against non-small cell lung cancer. Other approaches, on the other hand, fail
 298 to provide clear and human-understandable explanations.

299 5 Related Work

300 Local-level Graph Neural Network (GNN) explanation approaches have been developed to shed
 301 light on the decision-making process of GNN models at the individual data instance level. Most of
 302 them, such as GNNExplainer [31], PGExplainer [17], PGM-Explainer [27], GraphLime [10], RG-
 303 Explainer [23], CF-GNNExplainer [16], RCEExplainer [3], CF² [26], RelEx [34] and Gem [15], train a
 304 secondary model to identify crucial nodes, edges, or subgraphs that explain the behavior of pretrained
 305 GNNs for specific input samples. However, the quality of the explanations produced by these methods
 306 is highly dependent on hyperparameter choices. Moreover, these explainers' black-box nature raises
 307 doubts about their ability to provide comprehensive explanations for GNN models. Approaches like
 308 SA [4], Grad-CAM [20], GNN-LRP [22], and DEGREE [8], which rely on gradient back-propagation,
 309 encounter the saturation problem [24]. As a result, these methods may generate explanations that are
 310 less faithful. SubgraphX [33] combines perturbation-based techniques with pruning using Shapley
 311 values. While it can generate some high-quality subgraph explanations, its computational cost is
 312 significantly high due to the reliance on the MCTS (Monte Carlo Tree Search). Additionally, as
 313 demonstrated in our experiments in Section 4, existing approaches exhibit inconsistencies on similar
 314 data samples and poor discriminability. This reinforces the need for our proposed method *GOAt*,
 315 which outperforms state-of-the-art baselines on fidelity, discriminability and stability metrics. Our
 316 work also relates to global-level explainability approaches. GLGExplainer [2] leverages prototype
 317 learning and builds upon PGExplainer to obtain global explanations. GCFExplainer [11] generates
 318 global counterfactual explanations by employing random walks on an edit map of graphs, utilizing
 319 local explanations from RCEExplainer and CF². Both GLGExplainer and GCFExplainer heavily rely
 320 on local explanations. Integrating local explainers that produce higher-quality local explanations,
 321 such as *GOAt*, has the potential to enhance the performance of these global-level explainers.

322 6 Conclusion

323 We propose *GOAt*, a local-level GNN explainer that overcomes the limitations of existing GNN
 324 explainers, in terms of insufficient discriminability, inconsistency on same-class data samples, and
 325 overfitting to noise. We analytically expand GNN outputs for each class into a sum of scalar products
 326 and attribute each scalar product to each input feature. Although *GOAt* shares similar limitations
 327 with some decomposition methods of requiring expert knowledge to design corresponding explaining
 328 processes for various GNNs, our extensive experiments on both synthetic and real datasets, along
 329 with qualitative analysis, demonstrate its superior explanation ability. Our method contributes to
 330 enhancing the transparency of decision-making in various fields where GNNs are widely applied.

References

- 331
- 332 [1] Julia Amann, Alessandro Blasimme, Effy Vayena, Dietmar Frey, and Vince I Madai. Explainability for
333 artificial intelligence in healthcare: a multidisciplinary perspective. *BMC Medical Informatics and Decision*
334 *Making*, 20(1):1–9, 2020.
- 335 [2] Steve Azzolin, Antonio Longa, Pietro Barbiero, Pietro Lio, and Andrea Passerini. Global explainability of
336 GNNs via logic combination of learned concepts. In *The Eleventh International Conference on Learning*
337 *Representations*, 2023.
- 338 [3] Mohit Bajaj, Lingyang Chu, Zi Yu Xue, Jian Pei, Lanjun Wang, Peter Cho-Ho Lam, and Yong Zhang.
339 Robust counterfactual explanations on graph neural networks. *Advances in Neural Information Processing*
340 *Systems*, 34:5644–5655, 2021.
- 341 [4] Federico Baldassarre and Hossein Azizpour. Explainability techniques for graph convolutional networks.
342 In *International Conference on Machine Learning (ICML) Workshops, 2019 Workshop on Learning and*
343 *Reasoning with Graph-Structured Representations*, 2019.
- 344 [5] David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, and Antonio Torralba. Network dissection: Quantify-
345 ing interpretability of deep visual representations. In *Proceedings of the IEEE Conference on Computer*
346 *Vision and Pattern Recognition (CVPR)*, July 2017.
- 347 [6] Alexander Binder, Grégoire Montavon, Sebastian Lapuschkin, Klaus-Robert Müller, and Wojciech Samek.
348 Layer-wise relevance propagation for neural networks with local renormalization layers. In *Artificial*
349 *Neural Networks and Machine Learning–ICANN 2016: 25th International Conference on Artificial Neural*
350 *Networks, Barcelona, Spain, September 6-9, 2016, Proceedings, Part II* 25, pages 63–71. Springer, 2016.
- 351 [7] Lukas Faber, Amin K. Moghaddam, and Roger Wattenhofer. When comparing to ground truth is wrong: On
352 evaluating gnn explanation methods. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge*
353 *Discovery & Data Mining*, pages 332–341, 2021.
- 354 [8] Qizhang Feng, Ninghao Liu, Fan Yang, Ruixiang Tang, Mengnan Du, and Xia Hu. Degree: Decomposition
355 based explanation for graph neural networks. In *International Conference on Learning Representations*,
356 2022.
- 357 [9] Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs.
358 *Advances in neural information processing systems*, 30, 2017.
- 359 [10] Qiang Huang, Makoto Yamada, Yuan Tian, Dinesh Singh, and Yi Chang. Graphlime: Local interpretable
360 model explanations for graph neural networks. *IEEE Transactions on Knowledge and Data Engineering*,
361 2022.
- 362 [11] Zexi Huang, Mert Kosan, Sourav Medya, Sayan Ranu, and Ambuj Singh. Global counterfactual explainer
363 for graph neural networks. In *Proceedings of the Sixteenth ACM International Conference on Web Search*
364 *and Data Mining*, pages 141–149, 2023.
- 365 [12] Brian Kenji Iwana, Ryohei Kuroki, and Seiichi Uchida. Explaining convolutional neural networks using
366 softmax gradient layer-wise relevance propagation. In *2019 IEEE/CVF International Conference on*
367 *Computer Vision Workshop (ICCVW)*, pages 4176–4185, 2019.
- 368 [13] Jeroen Kazius, Ross McGuire, and Roberta Bursi. Derivation and validation of toxicophores for muta-
369 genicity prediction. *Journal of medicinal chemistry*, 48(1):312–320, 2005.
- 370 [14] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In
371 *International Conference on Learning Representations*, 2017.
- 372 [15] Wanyu Lin, Hao Lan, and Baochun Li. Generative causal explanations for graph neural networks. In
373 *International Conference on Machine Learning*, pages 6666–6679. PMLR, 2021.
- 374 [16] Ana Lucic, Maartje A Ter Hoeve, Gabriele Tolomei, Maarten De Rijke, and Fabrizio Silvestri. Cf-
375 gnnexplainer: Counterfactual explanations for graph neural networks. In *International Conference on*
376 *Artificial Intelligence and Statistics*, pages 4499–4511. PMLR, 2022.
- 377 [17] Dongsheng Luo, Wei Cheng, Dongkuan Xu, Wenchao Yu, Bo Zong, Haifeng Chen, and Xiang Zhang.
378 Parameterized explainer for graph neural network. *Advances in neural information processing systems*,
379 33:19620–19631, 2020.
- 380 [18] Xinjun Pei, Long Yu, and Shengwei Tian. Amalnet: A deep learning framework based on graph convolu-
381 tional networks for malware detection. *Computers & Security*, 93:101792, 2020.

- 382 [19] Douglas EV Pires, Tom L Blundell, and David B Ascher. pkcsm: predicting small-molecule pharmaco-
383 netic and toxicity properties using graph-based signatures. *Journal of medicinal chemistry*, 58(9):4066–
384 4072, 2015.
- 385 [20] Phillip E Pope, Soheil Kolouri, Mohammad Rostami, Charles E Martin, and Heiko Hoffmann. Explain-
386 ability methods for graph convolutional neural networks. In *Proceedings of the IEEE/CVF conference on*
387 *computer vision and pattern recognition*, pages 10772–10781, 2019.
- 388 [21] Michael Sejr Schlichtkrull, Nicola De Cao, and Ivan Titov. Interpreting graph neural networks for nlp with
389 differentiable edge masking. In *International Conference on Learning Representations*, 2021.
- 390 [22] Thomas Schnake, Oliver Eberle, Jonas Lederer, Shinichi Nakajima, Kristof T Schütt, Klaus-Robert Müller,
391 and Grégoire Montavon. Higher-order explanations of graph neural networks via relevant walks. *IEEE*
392 *transactions on pattern analysis and machine intelligence*, 44(11):7581–7596, 2021.
- 393 [23] Caihua Shan, Yifei Shen, Yao Zhang, Xiang Li, and Dongsheng Li. Reinforcement learning enhanced
394 explainer for graph neural networks. *Advances in Neural Information Processing Systems*, 34:22523–22533,
395 2021.
- 396 [24] Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. Learning important features through propagat-
397 ing activation differences. In *International conference on machine learning*, pages 3145–3153. PMLR,
398 2017.
- 399 [25] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *International*
400 *conference on machine learning*, pages 3319–3328. PMLR, 2017.
- 401 [26] Juntao Tan, Shijie Geng, Zuohui Fu, Yingqiang Ge, Shuyuan Xu, Yunqi Li, and Yongfeng Zhang. Learning
402 and evaluating graph neural network explanations based on counterfactual and factual reasoning. In
403 *Proceedings of the ACM Web Conference 2022*, pages 1018–1027, 2022.
- 404 [27] Minh Vu and My T Thai. Pgm-explainer: Probabilistic graphical model explanations for graph neural
405 networks. *Advances in neural information processing systems*, 33:12225–12235, 2020.
- 406 [28] Daixin Wang, Zhiqiang Zhang, Jun Zhou, Peng Cui, Jingli Fang, Quanhui Jia, Yanming Fang, and Yuan
407 Qi. Temporal-aware graph neural network for credit risk prediction. In *Proceedings of the 2021 SIAM*
408 *International Conference on Data Mining (SDM)*, pages 702–710. SIAM, 2021.
- 409 [29] Lingfei Wu, Peng Cui, Jian Pei, and Liang Zhao. *Graph Neural Networks: Foundations, Frontiers, and*
410 *Applications*. Springer Nature, 2022.
- 411 [30] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks?
412 In *International Conference on Learning Representations*, 2019.
- 413 [31] Zhitao Ying, Dylan Bourgeois, Jiaxuan You, Marinka Zitnik, and Jure Leskovec. Gnnexplainer: Generating
414 explanations for graph neural networks. *Advances in neural information processing systems*, 32, 2019.
- 415 [32] Hao Yuan, Haiyang Yu, Shurui Gui, and Shuiwang Ji. Explainability in graph neural networks: A
416 taxonomic survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.
- 417 [33] Hao Yuan, Haiyang Yu, Jie Wang, Kang Li, and Shuiwang Ji. On explainability of graph neural networks
418 via subgraph explorations. In *International Conference on Machine Learning*, pages 12241–12252. PMLR,
419 2021.
- 420 [34] Yue Zhang, David Defazio, and Arti Ramesh. Relax: A model-agnostic relational model explainer. In
421 *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, pages 1042–1049, 2021.