

GNN Interpretability Using Bayesian Inference

Shalin Patel

Advisor: Dr. Ritambhara Singh
Second Reader: Dr. Lorin Crawford



Division of Applied Mathematics and Department of Computer Science
Brown University
2023-04-18

Contents

1	Introduction	2
1.1	Graph Neural Networks	3
1.2	Interpretation on GNNs	4
1.3	Bayesian Inference	4
1.3.1	Stochastic Variational Inference (SVI)	6
1.3.2	Bayes-by-backprop	7
1.4	Normalizing Flows	7
2	Related Work	8
2.1	GNN Explainer	9
2.1.1	Benchmark Datasets	10
2.2	Parametrized-Graph Explainer	13
2.3	Other Explainer Frameworks	13
2.4	SERGIO	13
3	Methods	13
4	Experimental Setup	13
5	Results	13
6	Discussion	13

GNN Interpretability Using Bayesian Inference

Shalin Patel^{1,2}

¹*Division of Applied Mathematics, Brown University*

²*Department of Computer Science, Brown University*

April 4, 2023

Abstract

1 Introduction

Graphs serve as a natural repository for information in many real-world applications ranging from social, informational, chemical, and biological domains [1]. Especially as data becomes more and more unstructured, graphs represent a flexible manner for storing and relating different nodes and their related features [2]. Indeed, graphs represent one of the most general mathematical structures for relating data and are seeing increasing use in modeling phenomenon such as social networks and gene regulatory networks [2, 3]. For the purposes of this work, given a set of vertices V , node features $\mathcal{X} : V \rightarrow \mathbb{R}^d$, a set of edges $E \subseteq V \times V$, and weights on the edges $W : E \rightarrow [0, 1]$, we consider the graph $G = \{V, \mathcal{X}, E, W\}$. Additionally, we let the space of all graphs for a given set of vertices V be \mathcal{G} . Below in figure 1, a simple visualization of this definition can be seen for a cyclic graph of order 3.

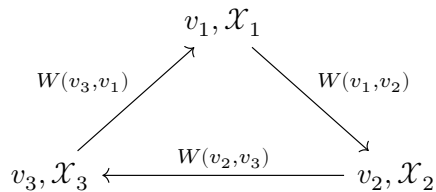


Figure 1: An example graph as related to the terminology laid out above

1.1 Graph Neural Networks

To deal with the proliferation of graphs in computing and the need to construct models that consider graphs as a first-class member of the modeling process, a class of models known as Graph Neural Networks have emerged (GNN) with state of the art performance on a variety of classification and regression tasks [4]. Specifically, GNNs and their early iterations in GCNs took inspiration from CNNs that represented applying successive convolution operations on regular grids of information, such as images, to compose local features in the grid in to higher-level predictions [5]. At a high level, most GNN frameworks can be split into three steps **MSG**, **AGG**, and **UPD** representing a messaging step, aggregation step, and update step, respectively.

At a layer l in a GNN model ϕ , the update of the hidden state of the model occurs first by sending messages for all $(v_i, v_j) \in E$ as a function of the hidden state \mathcal{H}_i^{l-1} and \mathcal{H}_j^{l-1} as well as the weight $W_{ij} := W(v_i, v_j)$. Specifically, we have

$$m_{ij}^l := \text{MSG}(\mathcal{H}_i^{l-1}, \mathcal{H}_j^{l-1}, W_{ij})$$

Then, a GNN performs an aggregation step wherein it calculates an aggregate message for every vertex $v \in V$. Let $\mathcal{N}_k : V \times E \rightarrow \mathcal{P}(E)$ be a function that returns the edges in the k -hop neighborhood of a node. Then, we can formally write the aggregation step as

$$M_i^l := \text{AGG}(\{m_{ij}^l \mid v_j \in \mathcal{N}_k(v_i)\})$$

Then, finally, at each node, the GNN takes a nonlinear function (often a neural network of some sort) and applies it to this aggregated message M_i^l along with the hidden state \mathcal{H}_i^{l-1} to get the new hidden state.

$$\mathcal{H}_i^l := \text{UPD}(\mathcal{H}_i^{l-1}, M_i^l)$$

When composed in layers, this forms a full Graph Neural Network. Note that in this framework, $\mathcal{H}_i^0 := \mathcal{X}_i$. Based on the task type, either node or graph classification in this work, further layers may be added on top of the final output. For example, in graph classification, it is often the case that the final node embeddings are concatenated and then run through an MLP to get a final classification for the whole graph [4].

The specific layers used in this work are Graph Convolution Layers also known as GCNs [6]. In the context of the framework above, the **MSG** sends weighted and normalized node features from the $l - 1$ st layer where the normalization is by the out-degree of the sending node and the weight of the message coming from W_{ij} . In the **AGG** step, all of these messages are summed together. Finally, the **UPD** step applies a neural network, usually a trainable linear layer combined with a non-linear activation function to provide a non-linear update step.

1.2 Interpretation on GNNs

Given a GNN, it is natural in many fields such as computational biology to perform interpretation on the model in order to gain further insights. For example, in the case of RNA-seq data, a natural question is to determine important regulatory pathways between genes that could be related to eventual up or down regulation of a target gene [3]. One natural way to perform this task is to train a GNN on the RNA-seq data for a node-classification task and feed it a large graph G with many redundant edges. Given a GNN model ϕ with l layers and a target gene $v_i \in V$, we would like to determine $\mathcal{E}_i \subseteq \mathcal{N}_l(v_i, E)$ as well as $\mathcal{W}_i : \mathcal{E}_i \rightarrow [0, 1]$ such that $\mathcal{E}_i, \mathcal{W}_i$ represent the most important subgraph for the model ϕ to perform its predictions for the input vertex v_i . While there are a few criteria for determining importance, we consider importance to be the maximal mutual information between the original model with its inputs and the same model with the estimated subgraph $\mathcal{E}_i, \mathcal{W}_i$. Written more formally, we wish to discover

$$\arg \max_{\mathcal{E}_i, \mathcal{W}_i} \text{MI}(\phi(v_i, \mathcal{X}, E, W), \phi(v_i, \mathcal{X}, \mathcal{E}_i, \mathcal{W}_i))$$

This is a computationally hard problem as a brute force search would take $O(2^{|E|})$ time even when discounting the weight array \mathcal{W}_i . In practice, discovering \mathcal{E}_i is ignored and most of the importance discovery is done through learning a suitable \mathcal{W}_i while letting $\mathcal{E}_i = E$.

However, because of the given task, and the various properties we would like to see in a reasonable interpretation of a GNN model, there are many routes that have been taken to interpret these models. As will be shown in §2, there have been a few non-bayesian attempts to solve this problem. The goal of this work is to analyze these methods and then suggest a new Bayesian method to solving the issue of searching for the best subgraph for a trained GNN to perform post-hoc analysis of importance. Furthermore, an added benefit of utilizing a Bayesian approach is that a full conditional distribution over the importance graph is learned giving researchers an even larger level of insight into their model that goes beyond just a simple edge mask.

1.3 Bayesian Inference

Recently, in deep learning literature there has been a rise in utilizing Bayesian methods to create Bayesian Neural Networks in order to provide uncertainty aware predictions [7]. While primarily concerned with giving estimates for failure modes and reducing overfitting within deep learning methods, the synthesis of Bayesian methods and deep learning models has imbued them with a greater sense of interpretability and introspectability. In the context of GNN interpretation, modeling the subgraph $\mathcal{E}_i, \mathcal{W}_i$ as a joint distribution allows for conditioning on certain edge weights and imbues the interpretations that are derived with a greater sense of introspectability. Figure 2 gives a good overview of the deep learning analogues for point estimate neural networks. In the same way, analogues will be utilized in the interpretation task in order to get the same benefits that have already been

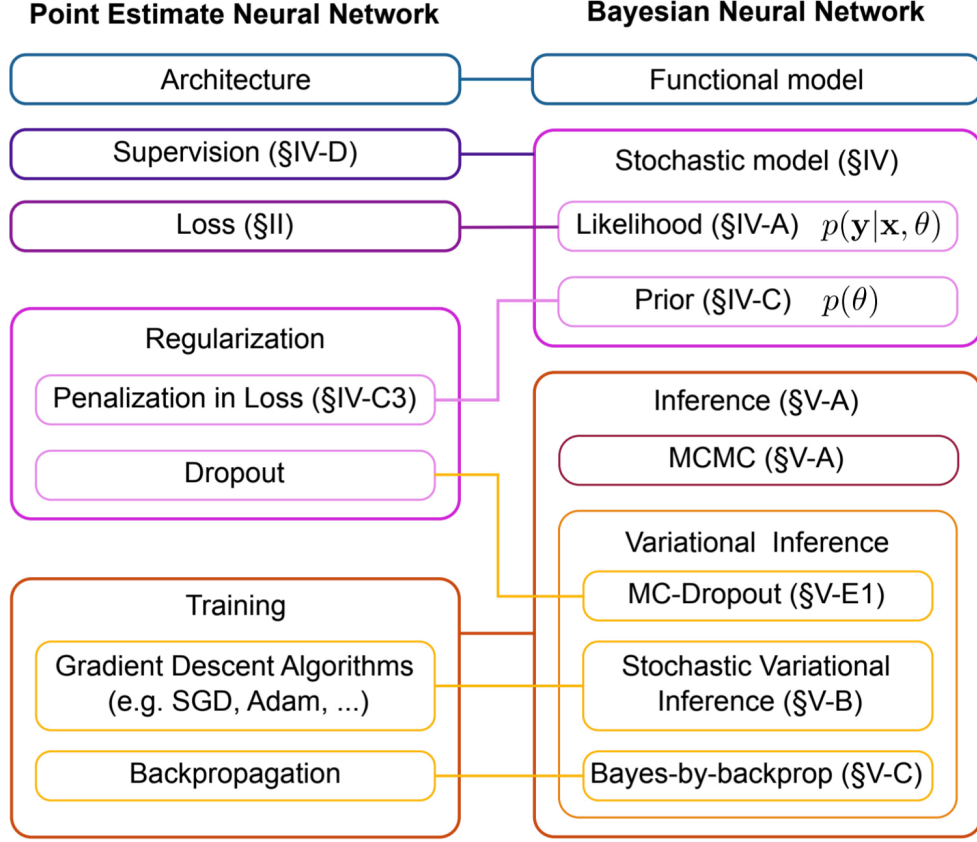


Figure 2: An overview of the corresponding structures in a standard neural network and a bayesian neural network

outlined in BNNs. In the Bayesian paradigm, a distribution \mathcal{P} is treated as the belief in the occurance of a given event from the distribution rather than the limit of the frequencies of each event as in the frequentist scheme. Furthermore, prior beliefs are thought to inform posterior beliefs. In the context of interprebility this is important as the belief for a given interpretation is dependent on the domain that it is brought up in. For example, in social networks one may expect relatively dense explanations while in biology they would tend to be sparse [1] [3]. Generally speaking, given a hypothesis H representing the prior belief for the state of a system and some data D , the posterior probability $\mathcal{P}(H | D)$ can be calculated as

$$\mathcal{P}(H | D) = \frac{\mathcal{P}(D | H)\mathcal{P}(H)}{\mathcal{P}(D)}$$

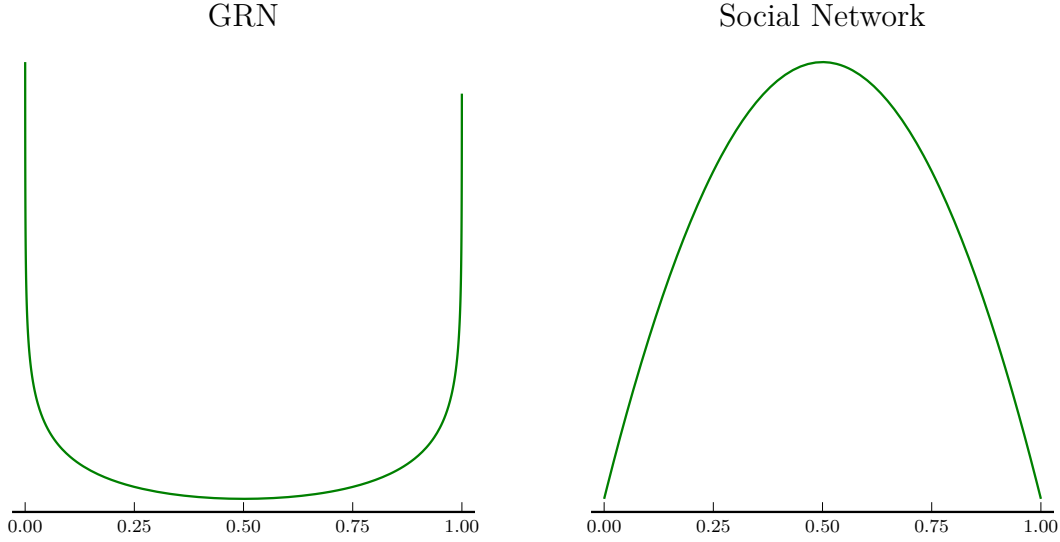


Figure 3: Idealized prior distributions for GRNs and Social Networks based on the Beta distribution. The displayed distributions are $B(0.95, 0.95)$ and $B(2, 2)$ respectively

and in this way, the posterior is conditioned by both the prior belief and the evidence that the data presents. While there are a variety of different techniques that can be used to update the prior belief into a posterior distribution as seen in figure 2, in this paper only stochastic variational inference (SVI) and Bayes-by-backprop are utilized.

1.3.1 Stochastic Variational Inference (SVI)

While other inference methods such as MCMC (with popular algorithms including HMC and NUTS [8] [9]), allow exact sampling from the posterior distribution, these methods have proven unpopular with the BNN community due to their algorithmic complexity and lack of scalability to larger models. Hence many communities use SVI which is not an exact method. In SVI, there is a family of distributions $q_\phi(H)$ which are parametrized by parameters ϕ . A common example would be the family of normal distributions parametrized by their mean and covariance structure. The goal of SVI is to approximate the posterior $\mathcal{P}(H | D)$ as closely as possible by $q_\phi(H)$. The most common measure of approximation in probability space is given by the Kullback-Leibler divergence (KL-divergence). While not a proper metric over the space of distributions, it does give a computationally-reasonable method to optimize against ϕ to get as close a match as possible. Specifically, SVI aims to minimize

$$D_{KL}(q_\phi || \mathcal{P}) = \int_H q_\phi(H') \log \frac{q_\phi(H')}{\mathcal{P}(H' | D)} dH'$$

This is still problematic since the quantity $\mathcal{P}(H \mid D)$ would still need to be calculated. Hence, it is sufficient to optimize against the ELBO which serves as a lower-bound for the KL-divergence. The ELBO is defined as

$$\log \mathcal{P}(D) - D_{KL}(q_\phi \parallel \mathcal{P}) = \int_H q_\phi(H') \log \frac{\mathcal{P}(H', D)}{q_\phi(H')} dH'$$

Note here that $\log \mathcal{P}(D)$ is just a constant meaning that minimizing the KL-divergence is the same as maximizing the ELBO. Note that, generally speaking, the families q_ϕ tend to come from the exponential family of distributions and the parameters for these families are then just optimized using a typical SGD algorithm such as ADAM [10].

1.3.2 Bayes-by-backprop

While SVI provides a good framework for Bayesian inference, it does not quite work for deep learning applications because stochasticity stops backpropagation from going through a neural network. To mitigate this problem, the usual reparametrization technique used in creating variational autoencoders (VAEs) [11] is combined with SVI to create a deep-learning friendly SVI algorithm. In this variation, a simple non-parametrized random variable $\epsilon \sim q(\epsilon)$ is sampled. To obtain the family $q_\phi(\theta)$, a deterministic transformation $t(\epsilon, \phi)$ is applied such that $\theta = t(\epsilon, \phi)$ has the property that $\theta \sim q_\phi(\theta)$. To obtain such a t , only a certain class of functions can be utilized. These functions are broadly known as bijectors and require t to be a diffeomorphism. In more detail, let $t : M \rightarrow N$ be a differentiable map, then t is a diffeomorphism if it is a bijection and its inverse $t^{-1} : N \rightarrow M$ is differentiable as well.

Generally speaking, the exponential family of distributions can all be constructed from such transformations meaning that they are good candidates for Bayes-by-backprop. Note though, that because of the transformation t , the formula for the ELBO changes to the following

$$\int_\epsilon q_\phi(t(\epsilon, \phi)) \log \frac{\mathcal{P}(t(\epsilon, \phi), D)}{q_\phi(t(\epsilon, \phi))} |\det(\nabla_\epsilon t(\epsilon, \phi))| d\epsilon$$

This is much friendlier to compute since ϵ is now a constant with respect to ϕ and lets us perform SVI through multiple layers of transformations simply by using bijectors like t .

1.4 Normalizing Flows

The technique that was described in §1.3.2 is more generally known as a normalizing flow. While it was described earlier in the context of a reparametrization technique in which the t are fixed, there is no such restriction in reality. More concretely, the t do not have to be simple functions but, rather, can be learnable functions in their own right. This allows one to use the normalizing flow technique to perform tasks like density estimation and

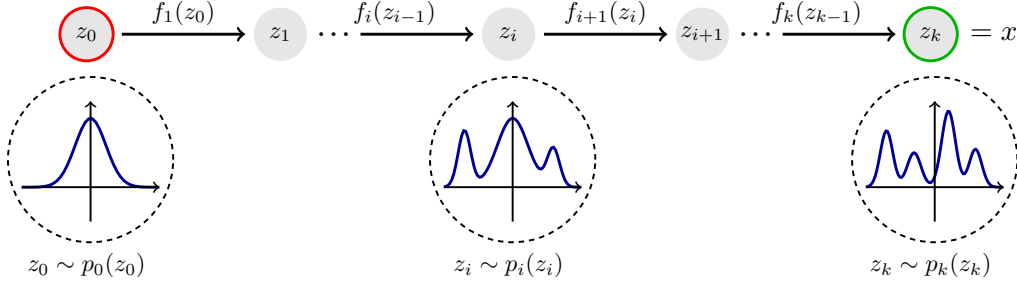


Figure 4: Chaining learnable bijectors to transform a base distribution to a more complicated distribution from [13]

distribution fitting with a very flexible class of transforms that take a simple distribution like a standard multivariate normal and make them into any computable distribution [12]. Let the transformations of the base distribution be defined as $g = g_n \circ g_{n-1} \circ \dots \circ g_1$ with inverse $f = g_1^{-1} \circ g_2^{-1} \circ \dots \circ g_n^{-1}$. Then we know that the determinant of the Jacobian of f is given by the product of the determinants of the Jacobians at each intermediate evaluation of the flow. This allows for more and more complicated transformations by introducing more and more learnable layers as can be seen in figure 4. This structure is very similar to that of an artificial feed-forward neural network. As an example, the simplest form of a normalizing flow is

$$g_i(x) = Ax + b$$

with the learnable parameters here being the matrix A and bias vector b . As long as A is an invertible matrix, we have a bijective function that can be used as a normalizing flow layer. Note that these linear layers can be interleaved with activation functions to provide non-linear transformations. This is important as a linear transformation of an exponential family will remain exponential so an element of non-linearity is required (as is the case with MLPs). While RELU is not invertible, a formulation like leaky-RELU can be used for this task [14]. Still these are not super expressive. For a normalizing flow with universality, this paper utilizes rational quadratic spline based flows [15]. When combined with variational inference over the prior base-distribution, this allows for a very flexible estimation of the posterior distribution no matter how complex.

2 Related Work

When it comes to GNN interpretability, there are a few main methods. The first that started the field of GNN interpretability was GNN Explainer [4] which also provided a suite of general benchmarks that most methods have utilized as a framework for analyzing the effectiveness of their explainer framework. Another major piece of work in the field

has been the parametrized-graph explainer (PGExplainer) [16] that took GNNExplainer and parametrized it with a deep neural network for faster inference times and more robust interpretation. Along with these two, a few other more recent explainers such as SubgraphX [17] and Gem [18] have introduced new ideas into the field with a variety of approaches to the problem of GNN interpretability. To date, there seems to be no fully Bayesian method to the problem of GNN interpretability.

In addition to these works, the work of SERGIO [19] will be introduced as it will be utilized later on to generate a new class of experiments that GNN Interpretability methods can be benchmarked against. This work provides causal graph structures that give a guaranteed groundtruth for interpretation.

2.1 GNN Explainer

The full version of GNNExplainer attempts to learn both a node interpretation and edge interpretation. For this work, only the edge interpretation part of the framework was utilized. GNN Explainer attempts to solve the objective outlined in §1.2, by only trying to learn \mathcal{W}_i while treating $\mathcal{E}_i = E$. In this framework, GNN Explainer enforces that for any $e \in E$, $W(e) \geq \mathcal{W}_i(e)$. Then to get the argmax, GNNExplainer treats \mathcal{W}_i as a random variable. Then the goal get transformed to

$$\arg \min_{\mathcal{W}_i} \mathbb{E}_{w \sim \mathcal{W}_i} [H(\phi(v_i, \mathcal{X}, E, w))]$$

This still remains intractable, so GNNExplainer attempts to make this simpler by using Jensen’s inequality. Note that this is not a reasonable application of Jensen’s since ϕ as a GNN has almost no hope of being convex. Nonetheless, using Jensen’s inequality gives

$$\arg \min_{\mathcal{W}_i} H(\phi(v_i, \mathcal{X}, E, \mathbb{E}[\mathcal{W}_i]))$$

This is still quite intractable if \mathcal{W}_i is a full joint distribution over all $e \in E$. Therefore, GNNExplainer attempts to use a mean field approximation for \mathcal{W}_i where the edge interpretation is decomposed into the product of Bernoulli distributions meaning that each edge weight is an independent Bernoulli distribution with mean equal to the underlying probability of the variable. Specifically,

$$\mathcal{P}(\mathcal{W}_i) = \prod_{(v_j, v_k) \in E} \mathcal{W}_i[v_j, v_k]$$

with each $\mathcal{W}_i[v_j, v_k]$ is a value between $[0, 1]$ representing a Bernoulli variable for each edge in the underlying graph as defined by E . In this case, if the classification for a node is c , GNNExplainer performs direct gradient descent on this array of values to minimize

$$\arg \min_{\mathcal{W}_i = \{\mathcal{W}_i[v_j, v_k] | (v_j, v_k) \in E\}} - \sum_{c=1}^C \mathbb{1}_{y=c} \log \mathcal{P}(\phi(v_i, \mathcal{X}, E, \mathcal{W}_i) = c)$$

While there is some probabilistic formulation here, in effect, GNNExplainer optimizes an adjacency matrix in $[0, 1]$ against the mutual information of the model given the edge weights and the model with the original graph. This means that GNNExplainer learns no conditional structure between edges and does not take the graph dynamics into account while training. This is further emphasized with the fact that a mean-field approximation was used to assume conditional independence between the edges of the underlying graph. Hence, it is an algorithm that provides only a summary of the interpretation using assumptions that are not generally applicable to GNNs.

2.1.1 Benchmark Datasets

As one of the first explainer methods for GNNs, GNNExplainer created a set of synthetic datasets that serve as the canonical datasets for GNN Interpretability. The goal of this section is to describe these datasets and the perceived shortcomings in these datasets that led to an exploration of their validity and a setup for the experiments produced later that demonstrate the incorrectness of these datasets for the stated task.

The main dataset focused on in this paper is the Tree-Cycles dataset [4]. In this dataset trees of depth three are attached to cycles of length six in order to form and aggregate dataset. A GNN node classification task entails predicting whether a given node is either in a tree portion of the graph or in the cyclic portion of the graph with no given node features. The idea here is that the GNN can only rely upon the structure of the graph for its node classification and all its information must come from the edges. Hence interpretation on the edges of the GNN would reveal only information that could be gleaned from the graph structure. In figure 5, one can see an example of a portion of the dataset looking at a three-hop neighborhood around node 565.

While this is a great construction to test an interpretation method, given that the GNN only relies on the graph structure for prediction, it is difficult to imagine what the ground truth for a given dataset is. The paper that introduced GNNExplainer proposed that the groundtruths be motifs in the graph. So if a node was in a cycle portion, the groundtruth would be the edges in the cycle and if a node was in the tree portion, the groundtruth would be the edges of the tree. This can be seen in figure 6 which shows the proposed ground truths. While this is a valid task for an interpretation technique to attempt to solve, it does not have direct bearing on the task that the GNN was trained on and it does not have direct bearing on the mutual information framework that forms the theoretical underpinning for GNN interpretability.

In a simple example, suppose a GNN is trained on the tree-cycles dataset. While it might be nice if the GNN needed all the edges in a cycle structure to determine that the node is, indeed, in the cycle, the GNN could have just as easily learned to use five edges from the cycle to make its predictions. Even if a theoretical GNN interpretation model was perfect, the fact that the GNN itself does not use the sixth edge means that an interpretation technique is doomed to max its accuracy score at $5/6$ which would make it

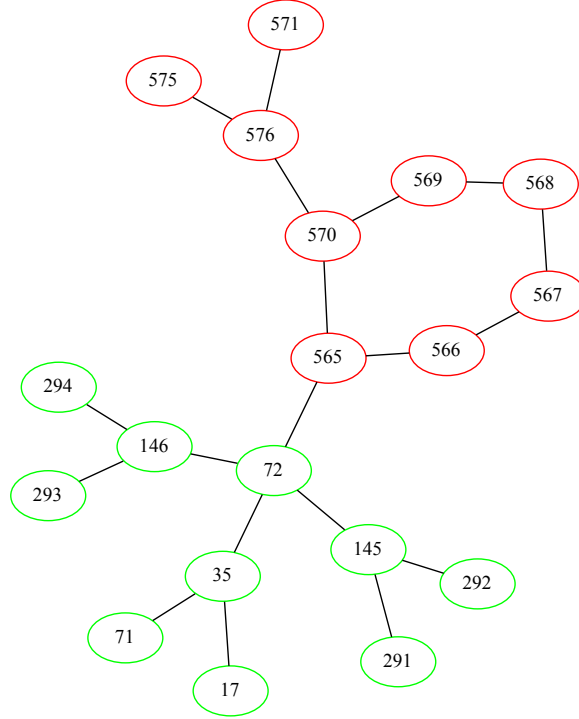


Figure 5: A look at the tree-cycles dataset at a three-hop neighborhood around node 565. In green are nodes classified as tree nodes and in red are nodes classified as cycle nodes. The dataset is almost 50% balanced between these two types.

impossible to compare between methods. Furthermore, a GNN may learn to mix motifs when making its prediction. Consider node 565 in the example from figure 5. While this node is in the cycle portion of the graph, the GNN could very easily learn that the node is *adjacent* to the tree structure detected in node 72 and learn to make the inverse decision in this case. Notably, there is no guarantee that a GNN learns the motifs as the important substructures and there is little chance that across nodes, it consistently learns these rules. Indeed, it will be shown later that GNNExplainer itself struggles to meet its $> 90\%$ accuracy scores for this benchmark in replication studies [17] [18] and that a thorough search through all connected subgraphs fails to yield this structure as the ground truth.

While GNNExplainer suggests a few other benchmark datasets, they all suffer from the same issue. Namely, they claim that the ground truth is the embedded motif structure, but a simple tests reveal that this is not consistently true and nor should it be true in the

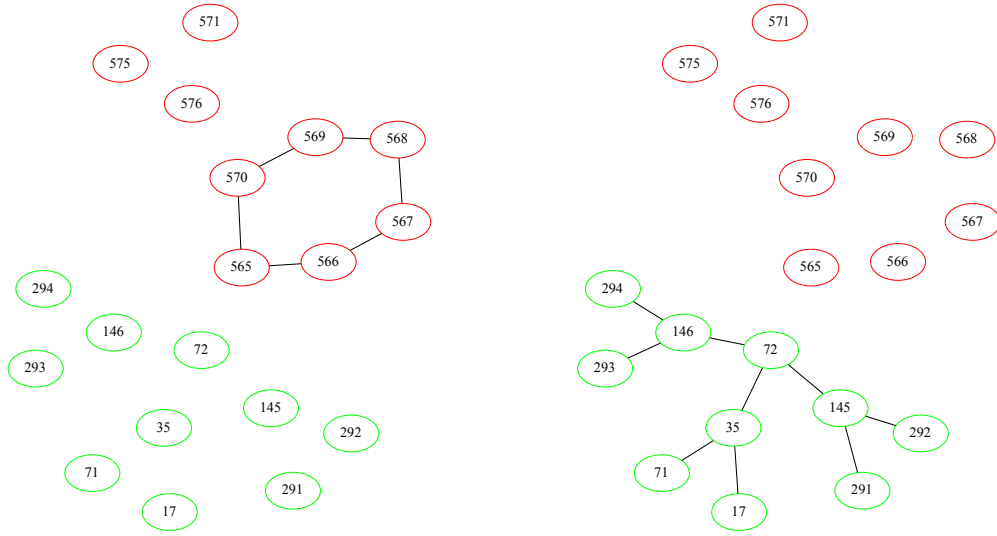


Figure 6: Demonstrations of the proposed groundtruths under [4] for the tree-cycles dataset. On the left is the proposed groundtruth for a node in the cycle (565) and on the right is the proposed ground truth for a node in a tree (72).

general case. Hence, the goal of this paper is to also provide a set of alternative benchmarks against which to evaluate GNN Interpretability.

2.2 Parametrized-Graph Explainer

2.3 Other Explainer Frameworks

2.4 SERGIO

3 Methods

4 Experimental Setup

5 Results

6 Discussion

References

- [1] Eunjoon Cho, Seth A. Myers, and Jure Leskovec. Friendship and mobility: user movement in location-based social networks. In *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*, KDD '11, pages 1082–1090. Association for Computing Machinery.
- [2] Takashi Washio and Hiroshi Motoda. State of the art of graph-based data mining. 5(1):59–68.
- [3] Francesca Petralia, Won-Min Song, Zhidong Tu, and Pei Wang. New method for joint network analysis reveals common and different coexpression patterns among genes and proteins in breast cancer. 15(3):743–754.
- [4] Rex Ying, Dylan Bourgeois, Jiaxuan You, Marinka Zitnik, and Jure Leskovec. GNNExplainer: Generating explanations for graph neural networks.
- [5] Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional neural networks on graphs with fast localized spectral filtering.
- [6] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks.
- [7] Laurent Valentin Jospin, Wray Buntine, Farid Boussaid, Hamid Laga, and Mohammed Ben-namoun. Hands-on bayesian neural networks – a tutorial for deep learning users. 17(2):29–48.
- [8] Michael Betancourt. A conceptual introduction to hamiltonian monte carlo.
- [9] Matthew D. Hoffman and Andrew Gelman. The no-u-turn sampler: Adaptively setting path lengths in hamiltonian monte carlo.
- [10] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization.
- [11] Diederik P. Kingma and Max Welling. An introduction to variational autoencoders. 12(4):307–392.

- [12] Ivan Kobyzev, Simon J. D. Prince, and Marcus A. Brubaker. Normalizing flows: An introduction and review of current methods. 43(11):3964–3979.
- [13] Lilian Weng. Flow-based deep generative models. Section: posts.
- [14] Bing Xu, Naiyan Wang, Tianqi Chen, and Mu Li. Empirical evaluation of rectified activations in convolutional network.
- [15] Conor Durkan, Artur Bekasov, Iain Murray, and George Papamakarios. Neural spline flows.
- [16] Dongsheng Luo, Wei Cheng, Dongkuan Xu, Wenchao Yu, Bo Zong, Haifeng Chen, and Xiang Zhang. Parameterized explainer for graph neural network.
- [17] Hao Yuan, Haiyang Yu, Jie Wang, Kang Li, and Shuiwang Ji. On explainability of graph neural networks via subgraph explorations.
- [18] Wanyu Lin, Hao Lan, and Baochun Li. Generative causal explanations for graph neural networks.
- [19] Payam Dibaeinia and Saurabh Sinha. SERGIO: A single-cell expression simulator guided by gene regulatory networks. 11(3):252–271.e11.