Explaining Graph Neural Networks with Mixed-Integer Programming

Anonymous Authors¹

Abstract

Graph Neural Networks (GNNs) provide stateof-the-art graph learning performance, but their lack of transparency hinders our ability to understand and trust them, ultimately limiting the areas where they can be applied. Many methods exist 015 to explain individual predictions made by GNNs, but there are fewer ways to gain more general insight into the patterns they have been trained to 018 identify. Most existing methods for model-level 019 GNN explanations attempt to generate graphs that 020 exemplify these patterns, but the discreteness of 021 graphs and the nonlinearity of deep GNNs make finding such graphs difficult. In this paper, we formulate the search for an explanatory graph as a mixed-integer programming (MIP) problem, in 025 which decision variables specify the explanation graph and the objective function represents the 027 quality of the graph as an explanation for a GNN's 028 predictions of an entire class in the dataset. This 029 approach, which we call MIPExplainer, allows us 030 to directly optimize over the discrete input space and find globally optimal solutions with a minimal number of hyperparameters. MIPExplainer outperforms existing methods in finding accurate 034 and consistent explanations on both synthetic and 035 real-world datasets.

1. Introduction

038

039

Graph neural networks (GNNs), such as graph convolutional networks (GCN) (Kipf & Welling, 2016), GraphSAGE net-041 works (Hamilton et al., 2017), and graph attention networks (GAT) (Veličković et al., 2017), provide a family of pow-043 erful tools for modelling graphs that learn from both the features contained in nodes and edges and the structure of 045 the graph itself. However, without being able to explain the 046 patterns GNNs rely on to make predictions, it is impossible 047 to justify their use in applications where trust and safety are 049 important, and there is no way to extract useful information from them. These problems have motivated a significant 050 051 amount of research into techniques for GNN explainability.

Research on explainable deep learning proceeds along two
 lines. One line is to develop intrinsically explainable meth-

ods, which modify standard neural networks or the training process so that final models naturally expose information about the importance and interaction of input features. Several proposed GNN architectures aim to achieve inherent explainability, for example, ProtGNN (Zhang et al., 2022) and GIB (Yu et al., 2020). The disadvantage of this approach is that changing the GNN itself to enforce explainability generally comes at the cost of performance. As a result, there is great interest in the second line of research, post-hoc explainability, which aims to interpret networks that have already been trained.

Post-hoc explanation for individual predictions has been extensively explored (see surveys from (Liu et al., 2021; Yuan et al., 2023; Kakkad et al., 2023)), but fewer methods exist to explain the overall patterns used by GNNs to differentiate classes (Yuan et al., 2020; Wang & Shen, 2022; Azzolin et al., 2022). There are several common problems among existing approaches for model-level GNN explanation, which focus on generating graphs to reflect knowledge learned by a model. For one, they often have many hyperparameters that can change the generated explanations, and generating a high-quality explanation may require setting them within a specific range of values. Without a single metric to quantify explanation quality, it is impossible to see if a certain choice of hyperparameters was effective, let alone compare the results across different hyperparameter settings. Furthermore, many methods randomly initialize parameters used to generate the explanations, and then rely on stochastic gradient optimization to assign them values. Due to the stochastic nature of the approach and the necessity of setting a maximum number of iterations, the final explanation's objective value might be far away from the global optimum. More importantly, it leads to significant variation in the explanations across different random initializations of learned parameters, even with the same choices of hyperparameters. Due to the lack of consistency and inability to guarantee solution quality, generating trustworthy global explanations in this way is almost impossible.

Because generating a graph is naturally a discrete process, we propose a new explanation method based on mixedinteger programs (MIPs), which we call MIPExplainer, for finding graph structures or subgraphs that maximally differentiate distinct classes as reported by the GNN. A MIP defines a constrained optimization problem where some of the

decision variables must take integer values. They are commonly solved through branch-and-bound, where the original 057 problem is split into subproblems that partition the set of 058 feasible solutions and solved recursively, creating a search 059 tree. Because an upper bound on any of these subproblems 060 can be found quickly by relaxing the integrality constraints, 061 large subtrees can be pruned if the upper bound at an in-062 ternal node is less than the objective value for a known 063 solution, improving the tractability of the search. Differing 064 from current methods that search for a graph maximizing 065 the output probability of a single class, we also define a 066 new explanation objective that measures the discriminative 067 power of the GNN. We further propose a new quantitative 068 metric to assess the consistency of the explanations from 069 multiple runs of an explanation method, by measuring the 070 dissimilarity of the generated graphs. While any appropriate graph distance metric can be used in conjunction with our framework (Gao et al., 2010), we employ graph edit distance (Sanfeliu & Fu, 1983b) which is commonly used 074 in inexact graph matching.

075 MIPExplainer offers several benefits over existing ap-076 proaches. (1) It directly optimizes over the discrete space 077 of possible input graphs, without any restrictions on types 078 of node and edge features. The only assumptions we make 079 about the space of graphs are bounds on the number of nodes and the magnitude of their features, and we do not 081 require any assumptions about the underlying distribution 082 of the training data. (2) It has a minimal number of hyper-083 parameters that influence the explanation (only the number of nodes of the explanation graph must be specified), facil-085 itating the application of our approach and mitigating the 086 effects of bias when analyzing the results. (3) We prove that 087 our MIP formulation has a globally optimal solution, and in many cases, we can find and verify this solution. In cases 089 where this is intractable, MIPExplainer can place an upper 090 bound on the optimal solution, guaranteeing the quality of 091 the generated explanation. 092

1.1. Related Work

093

094

095 GNN interpretation has been largely focused on instance-096 level explanation, which aims to explain the reasoning be-097 hind individual predictions. As identified in (Yuan et al., 098 2023), at least six categories of instance-level GNN ex-099 planation methods have been proposed so far: gradient-100 based (Pope et al., 2019), perturbation-based (Yuan et al., 2021; Luo et al., 2020; Ying et al., 2019; Schlichtkrull et al., 2020), surrogate (Vu & Thai, 2020), generation-based (Lin et al., 2021), decomposition (Schnake et al., 2021), and 104 counterfactual-based (Lucic et al., 2022) methods. These 105 methods do not immediately provide insights into the over-106 all patterns a GNN has identified, but it is possible to consolidate instance-level explanations to reveal model-level patterns. For example, we can employ purely statistical 109

methods to determine whether there are nodes/edges shared by a significant portion of the individual explanations. A more recent technique, GLGExplainer (Azzolin et al., 2022), finds smaller components of the extracted explanations that can be used to build logical expressions consistent with the overall GNN's predictions. However, these methods are limited by the scope of the training data, and can be influenced by bias in the dataset. By allowing for explanations to be out of distribution, we are able to isolate graph structures that may not appear by themselves without noise in the actual data. A direct model-level explanation can offer more faithful explanations, and is more useful for determining the degree of bias in the model itself.

Relatively few methods exist to explain GNNs at the model level. XGNN (Yuan et al., 2020) is the most widely used, and serves as the only baseline in several recent papers that focus on similar objectives (Azzolin et al., 2022; Saha et al., 2023a; Shin et al., 2022; Wang & Shen, 2022). These methods aim to generate graphs that exemplify graph structures used by a trained GNN for making classifications, without straying too far from the distribution of the training data where the model is not well-defined. XGNN trains a second neural network by reinforcement learning to generate graphs that obey explicit generation rules and maximize the original GNN's prediction for a specific class. GNNInterpreter (Wang & Shen, 2022) and GraphEx (Saha et al., 2023b) avoid training a second neural network by assuming that the graphs in the dataset are sampled from a set of underlying distributions parameterized by continuous latent parameters. In particular, GNNInterpreter defines an objective function similar to XGNN during training, maximizing a target class's logit while penalizing the distance between the embedding of the generated graph and the mean embedding of the training data to keep explanations in-distribution, and learns parameters through Monte Carlo gradient estimation. GCExplainer (Magister et al., 2021) and the work by (Xuanyuan et al., 2023) are also global explanation methods, but both focus on explaining GNNs via concept generation, which is a separate approach that relies on identifying patterns in activation maps of the training data.

In the past, discrete optimization techniques have been applied to deep neural networks, such as in (Cheon, 2022), (Botoeva et al., 2020), and (Ansari et al., 2022). However, these existing methods were only defined for layers consisting of a linear transformation with ReLU activations, and solve constraint satisfaction or optimization problems related to inverse design and verification. This work will reformulate the optimization problem for explainability and generalize the application of mixed-integer programming on standard neural networks to a range of GNN architectures.

2. MIPExplainer

111 Our model-level explanation seeks to optimize an input 112 graph G = (X, A) on which the GNN maximally differ-113 entiates one class from the rest, where X contains the 114 d attributes for each of N graph nodes as row vectors 115 and $A = (a_{ii})$ represents the N by N adjacency matrix. 116 We focus on the case of a binary adjacency matrix where 117 $a_{ij} \in \{0, 1\}$, so that $a_{ij} = 1$ indicates there is an edge be-118 tween nodes i and j. Let a GNN realize a function $f_c(G, \theta)$ 119 that maps G to the probabilities of several classes indexed 120 by c and θ contains all the learned parameters in the GNN. 121 During the GNN training, G is given and θ needs to be deter-122 mined, whereas in many model explanation methods, θ has 123 been fixed, and we optimize G (i.e., X and A) to maximize 124 $f_c(G, \theta)$ (or a related objective). 125

126 The proposed MIP will optimize G in terms of the values 127 of A and X. Each layer of the GNN imposes a set of con-128 straints in the MIP. We add decision variables to represent 129 the output of each layer and add constraints to represent 130 the computation in that layer. For example, for a fully con-131 nected layer, a new matrix of decision variables Y will be 132 added to the model and constrained with Y = WY' + b, 133 where Y' are the decision variables representing the outputs 134 of a previous layer and W and b are the model's learned 135 parameters in this layer. Since the outputs of subsequent lay-136 ers are constrained exactly, ultimately all of the constraints 137 define the feasible region of X and A. We place constraints 138 on nodes and edges (or entries of A) so that the derived 139 explanation forms a connected graph with valid features, 140 and we can further constrain X and A to reduce the number 141 of candidate solutions for a single graph since the GNN is 142 permutation equivariant with respect to the order of nodes. 143 In the subsequent sections, we provide a detailed derivation 144 of our MIP formulation by discussing the objective function 145 and the various constraints. 146

147148**2.1. Objective Function**

A typical objective function for explanations contains two 149 parts: a term related to class prediction and a regularizer that 150 enforces the generated explanations to be in-distribution. In 151 this paper, we decide not to apply any regularization in the 152 objective function in order to minimize the number of hy-153 perparameters; please see our detailed discussions in the 154 Appendix. While maximizing a single logit while disregard-155 ing the logits of other classes in the denominator (e.g., as 156 done by GNNInterpreter) is possible, this may lead to low 157 quality explanations in some circumstances. Predictions are made based on the difference between the logits, and 159 the absolute value of a single logit may be unrelated to the 160 prediction of the network. To illustrate this, after training 161 a GNN to classify star graphs and wheel graphs of varying 162 sizes (a task defined in (Wang & Shen, 2022)), we plotted 163 164

the logits it assigned to the training data for both classes in Figure 1. Note that the maximum logit for the wheel class is actually assigned to a correctly-classified star graph. Thus, simply maximizing the logit for wheels will not produce an effective explanation for the wheel class.

To accurately find classdiscriminative information, we should maximize the difference between the logit of the target class and the logit of the other classes. Maximizing the normalized probability, as done by XGNN, is possible but can lead to numerical instability due to improvements getting exponentially smaller as



Graphs in the Shapes Dataset

the magnitude of the logits increases. We can form an objective function as a linear combination of all logits but with a positive coefficient for only the target class, but an optimal solution may simply minimize one logit while leaving other logits close to or even greater than the logit of the target class, resulting in an incorrect explanation. To mitigate this problem, we can maximize the difference between the logit of the target class and the maximum of the other classes. In our observation, this approach is more effective, so we focus on discussing the following objective function:

$$\max_{G} \left(f_c(G,\theta) - \max_{i \neq c} (f_i(G,\theta)) \right), \tag{1}$$

where f_i denotes the *i*th output of the GNN before the application of the softmax function for classification.

2.2. Constraints

We make one crucial assumption about the node features, that their values are bounded by a constant M. We do not make assumptions on the node features or their distribution. We also require that the number of nodes in the explanation n is fixed in advance, which is the only hyperparameter that changes the optimization problem being solved.

From the range of existing GNN layers, we choose to focus first on GraphSAGE convolution layers, where the updated node representations X' after a layer are calculated from existing node representations X with the formula

$$X' = \sigma(XW_1 + \text{Aggregation}(X)W_2 + b).$$
 (2)

The aggregation on a node can be realized, for example, by summing its neighbors' feature vectors, i.e., Aggregation(X) = AX.

Assume that a GNN model has L_c GraphSAGE-based con-165 166 volution layers with sum aggregations and ReLU activa-167 tions, followed by a global feature-wise sum pooling layer 168 and L_f fully connected (FC) layers with ReLU activations. 169 In total, there are $L = L_c + 1 + L_f$ layers (indexed as $L_i, i \in 1, \ldots, L$). We will use the following notations: 170 the matrix of scalars, $W^{(i)}$, and the vector of scalars, $b^{(i)}$, 171 172 denote the GNN's matrix of learned weights and learned bias vector in layer *i*. For convenience, we also denote 173 $X^{(0)} = X$, where x_{ij} is the *j*th feature of node *i*. 174

175 We will also add the following decision variables to our 176 formulation and discuss how they are constrained shortly: 177 $\Phi^{(i)}$ represents the output of layer *i* before the activation 178 function, $X^{(i)}$ represents ReLU($\Phi^{(i)}$), the output of layer 179 *i*. We also represent $ReLU(-\Phi^{(i)})$ by $B^{(i)}$, while $Z^{(i)}$ are 180 binary indicators representing the truth value of $\Phi^{(i)} > 0$ 181 elementwise¹. The vector d (with some abuse of notation) 182 is an indicator representing whether each element of $\Phi^{(L)}$ 183 is the maximum element in the output of layer L except the 184 target class, i.e., dimension j of d is 1 when dimension j185 of $\Phi^{(L)}$ is the maximum, while the rest are all 0's. Here, y 186 represents the value of the maximum output of the GNN that 187 is not for the target class. Aside from the binary variables 188 A, $Z^{(i)}$, and d, all other variables are continuous. 189

To constrain $\Phi^{(i)}$ for the convolutional layers $(1 \le i \le L_c)$: 190

191

192

196

197 198

199

200

201

202

204

206

211

212

213

214

215 216

$$\Phi^{(i)} = X^{(i-1)}W_1^{(i)} + AX^{(i-1)}W_2^{(i)} + b^{(i)}.$$
 (3)

193 For the pooling layer $i = L_c + 1$ (with 1 representing a 194 vector of 1s): 195

$$\Phi^{(i)} = \mathbf{1}^T \Phi^{(i-1)},\tag{4}$$

and for the fully connected layers $(L_c + 1 < i \leq L)$:

$$\Phi^{(i)} = X^{(i-1)} W_1^{(i)} + b^{(i)}.$$
(5)

To constrain $X^{(i)}$ for all layers except the pooling and readout layers $(0 < i \leq L - 1, i \neq L_c + 1)$, we encode the ReLU output as follows:

$$X^{(i)} - B^{(i)} = \Phi^{(i)}, \tag{6}$$

$$X^{(i)} \le MZ^{(i)},\tag{7}$$

$$B^{(i)} \le M(1 - Z^{(i)}),\tag{8}$$

$$0 \le X^{(i)}, B^{(i)} \le M.$$
 (9)

208 For the pooling layer, we simply have that 209 $X^{(L_c+1)} = \Phi^{(L_c+1)}$. To constrain d_j and y: 210

$$y \ge X_{\neq c}^{(L)},\tag{10}$$

$$y \le X_{\neq c}^{(L)} + (\max(U_{X_{\neq c}^{(L)}})\mathbf{1} - L_{X_{\neq c}^{(L)}})(1-d), \quad (11)$$

$$\sum_{j} d_{j} = 1, d_{j} \in \{0, 1\},$$
(12)

217 ¹For elements of $\Phi^{(i)}$ exactly equal to 0, the corresponding 218 values of $Z^{(i)}$ can still be 0, but this will not affect the computation. 219

where $L_{X_{\neq c}^{(L)}}$ and $U_{X_{\neq c}^{(L)}}$ represent lower and upper bounds for the decision variables in $X^{(L)}$ excluding the output of class c. A method to calculate these bounds based on the bounds of the input will be discussed in a later section. Most of these constraints are linear in terms of the decision variables except Eq.(3) where decision variables A and $X^{(i)}$ multiply to form quadratic terms. Because A is binary, these terms can be equivalently reformulated into linear functions, which makes the optimization significantly easier. There are several ways to perform the linearization of quadratic terms with both continuous and binary variables. We describe one such method by change of variables (Kalvelagen, 2008). For a given binary variable $a \in A$ and a continuous variable $x \in$ $X^{(i)}$ bounded by M, let $e = a \times x$ be a new intermediate decision variable. Let $E^{(i)}$ be the matrix of $AX^{(i)}$ where entries are all calculated by summing the corresponding e's. Constraints in Eq.(3) can be rewritten as follows with additional bound constraints:

$$\Phi^{(i)} = X^{(i-1)}W_1^{(i)} + E^{(i)}W_2^{(i)} + b^{(i)}, \qquad (13)$$

$$-Ma \le e \le Ma,\tag{14}$$

$$x - M(1 - a) \le e \le x + M(1 - a).$$
 (15)

Now, our MIP has been transformed into a problem that maximizes Eq.(1), which can be calculated as $X_c^{(L)} - y$ subject to constraints Eq.(13-15) and Eq.(4-12). Note that this MIP has a convex objective function and all linear constraints when integrality is relaxed, so it is a mixed-integer linear program (MILP).

2.3. Additional Constraints on A and X

Additional constraints can be placed on A and X when generating explanations. For example, when the input space is actually graphs with one-hot features, we can constrain the sum of each row of $X^{(0)}$ to be equal to 1, and X can also be defined with binary or integer decision variables when appropriate. If the input graph is undirected, we can add the constraints $a_{ij} = a_{ji}$ for all i, j with $0 \le ij < n$ and i < j. We can prevent self-loops in the explanation by constraining the diagonal elements of A to be 0.

We also impose a partial ordering on the graph nodes $node_1, \ldots, node_n$ to ensure that the explanation graph is connected, or in the case of directed graphs, weakly connected (i.e. connected ignoring the directionality of the edges). We require that there is at least one edge between node_i and the set of nodes {node_i | i > j}. This becomes a constraint on A, and specifically for each i with $0 \le i < n$, we add the constraint $\sum_{\{j|i>j\}} (a_{ij} + a_{ji}) \ge 1$. These constraints also partially alleviate the effect of equivariance (where the same graph can have many different A), because they reduce the MIP's feasible region, but not the set of candidate graphs. This can easily be proven by showing that

for any (weakly) connected graph, the nodes can be ordered in a way that satisfies these constraints by running depthfirst-search (DFS) on such a graph ignoring edge directions. The *i*th node found by DFS must have been found from one of the 1 through i - 1 nodes, so this is always true.

2.4. Generalizing to more GNNs

226

242

243

244

258

259

261

262

263

264

265

266

267

268

269

270

271

272

273

274

227 Many highly performant GNN architectures can be perfectly 228 represented by linear and quadratic constraints, and many 229 more can be closely approximated. For example, if we 230 choose our aggregation function to be a feature-wise aver-231 age instead of a feature-wise sum, we can simply modify 232 constraint (4) as $\Phi^{(i)} = \mathbf{1}^T \Phi^{(i-1)} \frac{1}{N}$ for $i = L_c + 1$, . If mean aggregation is used in Eq. (2), we need another set of 233 234 decision variables $D^{(i)}$ for each layer, where row j of $D^{(i)}$ 235 will represent the feature-wise average of the neighbors of 236 node j. To properly constrain $D^{(i)}$, we add the constraint 237 $\mathbf{1}(\mathbf{1}^T A)D^{(i)} = AX$ to the model. The elements of D can 238 be distributed in the multiplications for the left hand side, 239 and then all the terms can be linearized as previously de-240 scribed. Now, constraint (3) can be changed to: 241

$$\Phi^{(i)} = X^{(i-1)}W_1^{(i)} + D^{(i)}W_2^{(i)} + \boldsymbol{b}^{(i)}$$

245 Consider a message passing layer from a Graph Isomor-246 phism Network (Xu et al., 2018), where updated node rep-247 resentations are calculated as $X' = h((A + (1 + \epsilon)I)X),$ 248 h is a neural network, and ϵ is a constant. We can split 249 this computation by constraining intermediate decision vari-250 ables according to the inner piece, $AX + ((1 + \epsilon)I)X$, and 251 the application of the neural network to those intermediate 252 variables, both of which we previously discussed how to ex-253 press with linear constraints. Additionally, piecewise-linear 254 approximations can be created for non-linear functions, al-255 lowing us to model different activation functions and the 256 convolutional layers in GCNs or GATs. 257

2.5. Existence of Global Optima

We now show that a globally optimal solution to the MIP problem described above always exists.

Theorem 2.1. Consider the MIP problem that maximizes the objective function $X_c^{(L)} - y$, with decision variables $A, X^{(i)}, \Phi^{(i)}, Z^{(i)}, B^{(i)}, E^{(i)}, d_j, y$ subject to constraints (4-15). A global optimum exists for this MIP.

Proof. Since all the constraints are linear equalities or inequalities (after linearizing the multiplication of binary and continuous variables in Eq. (3)), they define a polyhedral feasible region in the space parameterized by the decision variables if binary variables are relaxed to be $\in [0, 1]$. All decision variables in the MIP are bounded, either directly or in terms of bounds on the input variables A and X, so

the feasible set is a closed polytope (a compact set). The objective function is linear in terms of the decision variables as calculated by $X_c^{(L)} - y$. Therefore, the linear relaxation of the MIP must have an optimum on the compact feasible set. In addition, there are a finite number of integer solutions within the compact feasible region, so at least one of them must have the maximum objective value.

2.6. Optimization Algorithm

We can employ a branch and bound procedure, along with cutting planes and heuristics, to find a globally optimal solution efficiently. We present a branch and bound algorithm in Algorithm 1 to solve our MIP described in the previous section, which is represented by its set of constraints C and objective function f. We obtain an initial solution at the root of a search tree by choosing an initialization graph $G_0 = (X_0, A_0)$ and applying the GNN to obtain initial values for all the intermediate variables.

Algorithm 1 MIP Branch and Bound Pr	rocedure
-------------------------------------	----------

- 1: **Input:** The constraint set C, the objective function f, and an initial graph $G_0 = (X_0, A_0)$
- 2: Initialize a queue Q containing only C
- 3: $L \leftarrow f(G_0)$
- 4: $z \leftarrow G_0$
- 5: while Q is not empty do
- 6: $N \leftarrow$ search node popped from Q
- 7: Solve the linear relaxation of N, denoted N_r , and store the result in z^*
- 8: $U \leftarrow f(z^*)$
- 9: **if** N_r was feasible and U > L **then**
- 10: **if** z^* obeys all integrality constraints, defining a valid graph G^* **then**
 - $z \leftarrow z^*$
 - $L \leftarrow f(G^*)$
- 13: else

11:

12:

14: $v \leftarrow An$ integer variable with a non-integral value z_v^* in z^*

15: Add $N \cup \{v \le \lfloor z_v^* \rfloor\}$ and $N \cup \{v \ge \lceil z_v^* \rceil\}$ to Q

```
16: end if
```

```
17: else
```

- 18: Prune the entire subtree rooted at N by continuing to the next iteration without adding any nodes to Q
- 19: end if

```
20: end while
```

```
21: return z
```

We start by finding the optimal solution of the LP relaxation (i.e. the MIP with the integrality constraints removed) of the MIP problem, for example using the simplex method (line 7). This is an optimal solution to a problem with fewer constraints, so it serves as an upper bound to the original
problem with integer domains for some decision variables.
If it happens to also be a solution to the original MIP, meaning all the variables are integral, then the correctness of the
algorithm for solving the relaxation guarantees that this is
an optimal solution to our MIP, and we can stop searching
(lines 10-12).

282 If some integral decision variables take fractional values in 283 the relaxation, we branch on one of them by partitioning 284 the set of candidate solutions for the problem with integer 285 constraints into 2 subproblems, one with the extra constraint 286 that a chosen fractional variable is at most the floor of its 287 value in the LP relaxation's optimum and another ensuring 288 that the variable at least the ceiling of that value (lines 14-289 15). The optimal solution to the integral problem will be 290 the maximum optimal solution of these two subproblems, 291 which can be solved recursively in the same way, leading to 292 a binary tree in which nodes represent further constrained 293 versions of the original MIP. While an upper bound for an 294 internal node (U in the pseudocode) only applies to the 295 subtree rooted at that node, any integer solution serves as a 296 lower bound (L in the pseudocode) to the integral problem's 297 optimum anywhere in the search tree. 298

299 Integer solutions can be found in the leaves of our search 300 tree, or using heuristics to complete partial solutions de-301 fined at internal nodes. If the linear relaxation solved at 302 an internal node is infeasible or has a maximum objective 303 value that is lower than or equal to our current lower bound, 304 we have proven that a new optimal solution to the integral 305 problem cannot lie anywhere in the subtree rooted at that 306 node, allowing us to prune the branch and skip all the nodes 307 it contains in our search (line 18). The process stops when 308 there are no more subproblems to explore, at which point we 309 will have found an optimal solution to the integral problem. 310 In our experiments, we use Gurobi Optimizer (Gurobi Opti-311 mization, LLC, 2023) to find an optimal solution efficiently. 312 This solver combines branch-and-bound with cutting plane 313 methods, which makes the optimization process even faster. 314 While the theoretical complexity of this algorithm is ex-315 ponential, the average complexity is significantly lower in 316 practice, making it tractable to apply in many situations.

317 A single, large number M can be used to bound all of the 318 continuous decision variables, but tighter bounds greatly 319 reduce the time needed to compute optimal solutions. While 320 automated bound-tightening procedures exist, it is faster to 321 use knowledge of the problem to bound manually. Each 322 hidden representation computed by the model is encoded 323 by a separate set of decision variables. Assuming we have 324 bounded the decision variables for one, we can compute 325 bounds for the outputs of a following transformation. For example, given a hidden representation vector x with lower 327 bound x_L and upper bound x_U , we can get upper and lower 328

329

bounds on the output of a linear layer x' = Wx + b:

$$x'_{L} = \operatorname{ReLU}(W)x_{L} + \operatorname{ReLU}(-W)x_{U} + b; \qquad (16)$$

$$x'_U = \operatorname{ReLU}(W)x_U + \operatorname{ReLU}(-W)x_L + b.$$
(17)

Given bounds on the the decision variables representing the explanation graph, the input to the GNN, we can follow the propagation of values through the GNN to iteratively bound the set of decision variables for each hidden representation. Bounds for the outputs of ReLU activation layers will be the same as their inputs, but clipped below at 0. In the case of layers like GraphSAGE convolutions where the output is the sum of several matrix multiplications, bounds can be derived for each term in the sum and then added together.

Further discussion on the practical considerations of solving these MIPs can be found in Appendix Section C.

3. Experiments

We use two synthetic datasets and one real-world dataset to evaluate our method: Is_Acyclic, Shapes, and MUTAG. The Is_Acyclic dataset comes from XGNN's experiments, and has two classes consisting of cyclic and acyclic graphs of various types. The cyclic graphs include graphs like grids, single cycles, and wheels, while the acyclic class includes graphs like paths and various types of trees. Every node is given the same feature, a single constant, in order to isolate the explanation methods' ability to capture structural information. For the Shapes dataset, which comes from GN-NInterpreter's experiments, graphs are first generated from one of 5 base classes: lollipop graphs contain a fully connected component with one connection to a path graph's end node, grid graphs are lattices where each internal node has 4 neighbors, star graphs have multiple outer nodes connected to a single central node, and wheel graphs are star graphs with a single cycle connecting the outer nodes. For each of these graphs, a uniform proportion between 0 and 0.2 is chosen, and the number of edges in the graph is increased by that amount by adding in edges uniformly at random. The features of each node are the same as in Is_Acyclic. The MUTAG dataset (Debnath et al., 1991) consists of graphs of chemical compounds, where nodes represent atoms and edges represent bonds between them. Each compound is classified as being either mutagenic or non-mutagenic. As described by the creators of this dataset and in (Hsu et al., 2016), mutagenic molecules tend to have higher numbers of fused rings of carbon atoms. For this dataset, each node's features are a one-hot vector indicating atom type.

To quantify the variation between explanation graphs, we run repeated experiments with each explanation method and measure the average graph edit distance between all pairs of explanations. Graph edit distance, as described in (Sanfeliu & Fu, 1983a), is the minimum number of graph edit operations (vertex/edge insertions/deletions/substitutions)

Submission and Formatting Instructions for IC	ML	2024
---	-----------	------

	# of Graphs	# of Classes	Average # of Nodes	Average # of Edges	# of N Feat
Is_Acyclic	533	2	28.5	68.1	
Shapes	8000	5	27.2	144.9	
	т	-1.1. 1 F	N-44 C		
	Ta	able 1. D	ataset Sumn	nary	
	Ta Train Ac	able 1. D	Dataset Sumn Test Accuracy	nary # of Model	Paramet
Is_Acyclic	Ta Train Ac	able 1. D	Dataset Sumn Test Accuracy 1.000	nary # of Model	Paramete 7
Is_Acyclic Shapes	Tain Ac	able 1. D ccuracy 0.998 0.991	Dataset Sumn Test Accuracy 1.000 0.993	nary # of Model	Paramet

Table 2. Performance Metrics of Trained GNNs

needed to transform one graph into another. A lower average graph edit distance indicates a more consistent explanation method.

3.1. Experimental Setup 345

Every dataset was randomly split into a training set (80%) 346 and a test set (20%). GNNs were trained on each, and perfor-347 mance metrics are reported in Table 2. For the experiments 348 with XGNN, we used the implementation provided by the 349 authors in DIG² (Liu et al., 2021). For the experiments with 350 GNNInterpreter, we also use the implementation provided 351 by the authors³. We used default sets of hyperparameters 352 provided in the papers and implementations of the baselines. 353 An exception was made for XGNN because the default reg-354 ularization weights sometimes caused the graph generator 355 to quickly learn a policy that stopped after the first node in several instances. To fix this, we increased the reward 357 for creating additional valid edges until it became favorable 358 for the model to generate reasonably sized explanations. 359 Additional details about the experiments can be found in 360 Appendix B. To measure the consistency of the explanations, 361 we generated explanations with 5,6,7, and 8 nodes using 362 each method on each dataset 5 times. Then, we computed 363 the average GED among the 5 explanations. Table 4 shows these metrics averaged over the different numbers of nodes. A full table containing separate results for each number of nodes can be found in Appendix D. For MIPExplainer, 367 the initialization graph was created by starting with a path graph and adding all other possible edges with probability 369 0.5. Variation in the baselines comes from the random ini-370 tialization of the explanation network for XGNN and the 371 latent parameters in GNNInterpreter.

3.2. Results

375 The main results from our experiments are shown in Table 3. 376 Note that when depicting molecular graphs, the node colors 377 are assigned as follows: gray=C, blue=N, red=O, cyan=F, 378 purple=I, green=Cl, and brown=Br. In the experiments with 379 Is_Acyclic, MIPExplainer explains the cyclic class with a 380

381

373

374

complete graph, which has the maximum possible number of cycles. It explains the acyclic class with a star graph, which is one of the most straightforward examples from the class. In contrast, the explanation graphs of XGNN and GNNInterpreter for the cyclic class contain some nodes with a single neighbor, and their explanations for the acyclic class even include multiple cycles. Our solver was able to prove the optimality of both classes, taking an average of 1.90 seconds for the cyclic class explanation and 119.14 seconds for the acyclic class explanation. The left plot in Figure 5 shows how the bounds converged over the course of the acyclic class experiment. This demonstrates how graph symmetries factor into MIPExplainer's runtime, we further discuss this problem in Appendix C. A fully connected graph with equal node features only has a single adjacency matrix and feature matrix representation, while a star graph with n nodes has n representations, as there are n options for the position of the central node in the node ordering. As a result, despite the solution having the same number of nodes and fewer edges, more of the search tree must be explored to prove the optimality of the acyclic explanation. Figure 5 shows the convergence of the objective bounds and the number of explored and unexplored search nodes over the course of the search for the optimal solution. Similar figures for more experimental settings can be found in Appendix D.

For the mutagenic class of the MUTAG dataset, the MIPExplainer produces a complete graph of carbon atoms. While the presence of carbon cycles are an important factor in the mutagenicity of organic molecules, they appear exclusively as rings of 5 or 6 carbon atoms. Neither of the explanations generated by the two baseline methods contained a cycle of carbon atoms. The explanations of the non-mutagenic class are not as reasonable across all methods, which is expected since non-mutagens are more accurately described by the absence of mutagenic features than by the presence of non-mutagenic features. The generated explanation mostly consists of bromine atoms, which only actually appear in 2 of the graphs in the dataset. Despite the larger network architecture, both these solutions were able to be verified as optimal within the time limit.

For the Shapes dataset, we can easily recognize the classes of each of the explanations generated by MIPExplainer. Despite the fact that a significant amount of noise was added to the training data, the explanations are relatively clean. For reference, three examples from the wheel class are shown in Figure 2. In the cases of lollipops and stars, we see important features of the graph duplicated, a lollipop with two ends and star with two centers. Because we chose a higher number of nodes to make the patterns clearer, optimality was not proven for these explanations, but they still appear reasonable.

Table 4 shows that MIPExplainer was significantly more

33 33

33

340

341

342

²https://github.com/divelab/DIG

³⁸² ³https://github.com/yolandalalala/

GNNInterpreter 383

³⁸⁴



Table 3. Generated Explanations. The graphs from left to right are generated by MIPExplainer, XGNN, and GNNInterpreter, respectively



Figure 2. Three randomly selected wheel graphs from the Shapes

399

400

401

402

403

404

405

406

407

408

409

435

436



410 Runtime (s) Runtime (s)
411
412 Figure 3. Solver metrics for 3 runs explaining the Acylcic class of
413 Is_Acyclic with 7 nodes: On the left, the current best solution's
414 objective (blue) and upper bound (red) converging to the same
415 global optimum (the dotted black line). On the right, the number

of explored (green)/unexplored (orange) nodes during the search. 416 417 consistent than the baseline methods in all experimental 418 settings. In some cases, small variations in explanations are 419 due to the existence of multiple explanation graphs with the 420 exact same objective value, which tended to be extremely 421 similar. In other cases, it was due to the algorithm running 422 out of time before finding the optimal solution. However, 423 this also rarely caused deviations, as the best solution was 424 generally found much earlier than it was proven to be op-425 timal. For example, proving optimality for the MUTAG 426 dataset for an explanation graph with 8 nodes always took 427 longer than 30 minutes, but across all 5 runs with random 428 initializations, the algorithm consistently produced the same 429 explanation after 30 minutes. An extended table showing 430 separate consistency metrics for differently sized explana-431 tions, as well as tables showing the output logits for the 432 explanation graphs and runtimes for each experiment, can 433 be found in Appendix D. 434

4. Conclusion and Discussion

437 Despite the ability of GNNs to model complex patterns in
438 graph-structured data, their lack of transparency remains one
439

		Averaged Edit Distance MIPExplainer GNNInterpreter XGNN			
Dataset	Class		•		
Is_Acyclic	Acyclic Cyclic	$\begin{array}{c}\textbf{0.2 \pm 0.40}\\\textbf{0.0 \pm 0.00}\end{array}$	$\begin{array}{c} 3.5 \pm 0.77 \\ 3.1 \pm 1.05 \end{array}$	$\begin{array}{c} 2.8 \pm 1.62 \\ 3.2 \pm 1.82 \end{array}$	
MUTAG	Mutagen Nonmutagen	$\begin{array}{c}\textbf{0.7}\pm\textbf{1.40}\\\textbf{0.4}\pm\textbf{0.80}\end{array}$	$\begin{array}{c} 8.1 \pm 1.19 \\ 7.8 \pm 1.56 \end{array}$	$\begin{array}{c} 7.8 \pm 2.17 \\ 7.4 \pm 2.39 \end{array}$	
Shapes	Grid Lollipop Star Wheel	$\begin{array}{c} \textbf{0.0} \pm \textbf{0.00} \\ \textbf{0.2} \pm \textbf{0.40} \\ \textbf{0.0} \pm \textbf{0.00} \\ \textbf{0.75} \pm \textbf{0.90} \end{array}$	$\begin{array}{c} 3.6 \pm 0.89 \\ 3.5 \pm 0.77 \\ 3.9 \pm 1.21 \\ 3.3 \pm 1.04 \end{array}$	$\begin{array}{c} 3.0 \pm 2.42 \\ 3.4 \pm 2.03 \\ 3.0 \pm 2.60 \\ 4.0 \pm 3.32 \end{array}$	

Table 4. Average edit distance between 5 generated example graphs, averaged for numbers of nodes between 5 and 8 inclusive. Time limit is one half hour.

of the key factors hindering their application in a wide range of domains. Model-level explanations of these networks are key to understanding the information they learn and improving their trust and reliability. In order to address key shortcomings that limit the use of existing methods in most real-world situations, this work proposes MIPExplainer for generating such explanations. Without a way to objectively evaluate their quality, it is essential that generated explanations are truly high-quality solutions of optimization problems that are not sensitive to user-defined hyperparameters. MIPExplainer achieves this by avoiding the use of both weighted regularizers and stochastic optimization, instead focusing on maximizing a simpler objective with deterministic methods that are able to prove the global optimality of the generated solutions. Minimal assumptions are made about the distributions of graphs and their features, and no secondary models are trained in the process.

The proposed method also has several shortcomings, which we hope to address in future work. While it is more general than previous methods in some ways, it also requires different GNN layers to be individually encoded with constraints, and may require piecewise-linear approximations for highly nonlinear components. From a practical perspective, the runtime of MIPExplainer as described here is the most significant drawback. Reducing symmetries in the encoding can greatly improve runtime, but this is a hard problem in general, and more work is required to understand which symmetries are the most costly when optimizing over sets of graphs. Despite these limitations, even before proving optimality, we observe that the proposed method is able to find reasonable explanations.

440 Impact Statement

441

442

443

444

445

446

447

448

449 450

451

452

453

454

455

456

457

478

Techniques for improving the explainability of deep neural networks have significant implications for safety and transparency wherever this technology is applied. While the method proposed in this work may be useful for elucidating implicitly learned patterns and biases, it is not able to make any guarantees about fairness or robustness, and should not be included in the verification process necessary for the safe deployment of deep learning models.

References

- Ansari, N., Seidel, H.-P., and Babaei, V. Mixed integer neural inverse design. ACM Transactions on Graphics, 41(4):151:1–151:14, July 2022. ISSN 0730-0301. doi: 10.1145/3528223.3530083. URL https://dl.acm. org/doi/10.1145/3528223.3530083.
- Azzolin, S., Longa, A., Barbiero, P., Liò, P., and Passerini,
 A. Global Explainability of GNNs via Logic Combination of Learned Concepts. 2022. doi: 10.48550/
 ARXIV.2210.07147. URL https://arxiv.org/ abs/2210.07147.
- Botoeva, E., Kouvaros, P., Kronqvist, J., Lomuscio, A., and Misener, R. Efficient Verification of ReLU-Based Neural Networks via Dependency Analysis. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(04):3291– 3299, April 2020. ISSN 2374-3468, 2159-5399. doi: 10. 1609/aaai.v34i04.5729. URL https://ojs.aaai. org/index.php/AAAI/article/view/5729.
- 472 Cheon, M.-S. An outer-approximation guided optimiza473 tion approach for constrained neural network inverse
 474 problems. *Mathematical Programming: Series A and*475 *B*, 196(1-2):173–202, November 2022. ISSN 0025476 5610. doi: 10.1007/s10107-021-01653-y. URL https:
 477 //doi.org/10.1007/s10107-021-01653-y.
- Debnath, A. K., Lopez De Compadre, R. L., Debnath, 479 G., Shusterman, A. J., and Hansch, C. Structure-480 activity relationship of mutagenic aromatic and het-481 eroaromatic nitro compounds. Correlation with molec-482 ular orbital energies and hydrophobicity. Journal 483 of Medicinal Chemistry, 34(2):786-797, February 484 1991. ISSN 0022-2623, 1520-4804. doi: 10.1021/ 485 URL https://pubs.acs.org/ jm00106a046. 486 doi/abs/10.1021/jm00106a046. 487
- Fey, M. and Lenssen, J. E. Fast Graph Representation Learning with PyTorch Geometric. In *ICLR Workshop on Representation Learning on Graphs and Manifolds*, 2019.
- 493 494 Gao, X., Xiao, B., Tao, D., and Li, X. A survey of graph edit

distance. *Pattern Analysis and applications*, 13:113–129, 2010.

- Gurobi Optimization, LLC. Gurobi Optimizer Reference Manual, 2023. URL https://www.gurobi.com.
- Hamilton, W. L., Ying, R., and Leskovec, J. Inductive Representation Learning on Large Graphs. June 2017. doi: 10.48550/arXiv.1706.02216. URL https: //arxiv.org/abs/1706.02216v4.
- Hsu, K.-H., Su, B.-H., Tu, Y.-S., Lin, O. A., and Tseng, Y. J. Mutagenicity in a Molecule: Identification of Core Structural Features of Mutagenicity Using a Scaffold Analysis. *PLOS ONE*, 11(2):e0148900, February 2016. ISSN 1932-6203. doi: 10.1371/journal. pone.0148900. URL https://dx.plos.org/10. 1371/journal.pone.0148900.
- Kakkad, J., Jannu, J., Sharma, K., Aggarwal, C., and Medya, S. A Survey on Explainability of Graph Neural Networks, June 2023. URL http://arxiv.org/abs/2306. 01958. arXiv:2306.01958 [cs].
- Kalvelagen, E. Multiplication of a continuous and a binary variable, May 2008. URL http: //yetanothermathprogrammingconsultant. blogspot.com/2008/05/ multiplication-of-continuous-and-binary. html.
- Kingma, D. P. and Ba, J. Adam: A Method for Stochastic Optimization. In Bengio, Y. and LeCun, Y. (eds.), 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings, 2015. URL http: //arxiv.org/abs/1412.6980.
- Kipf, T. N. and Welling, M. Semi-Supervised Classification with Graph Convolutional Networks. November 2016. URL https://openreview.net/forum? id=SJU4ayYgl.
- Lin, W., Lan, H., and Li, B. Generative causal explanations for graph neural networks. In *International Conference on Machine Learning*, pp. 6666–6679. PMLR, 2021.
- Liu, M., Luo, Y., Wang, L., Xie, Y., Yuan, H., Gui, S., Xu, Z., Yu, H., Zhang, J., Liu, Y., Yan, K., Oztekin, B., Liu, H., Zhang, X., Fu, C., and Ji, S. DIG: A Turnkey Library for Diving into Graph Deep Learning Research. arXiv preprint arXiv:2103.12608, 2021.
- Lucic, A., Ter Hoeve, M. A., Tolomei, G., De Rijke, M., and Silvestri, F. Cf-gnnexplainer: Counterfactual explanations for graph neural networks. In *International Conference on Artificial Intelligence and Statistics*, pp. 4499–4511. PMLR, 2022.

Luo, D., Cheng, W., Xu, D., Yu, W., Zong, B., Chen, H., 495 496 and Zhang, X. Parameterized explainer for graph neural 497 network. Advances in neural information processing 498 systems, 33:19620-19631, 2020.

499

506

527

528

- 500 Magister, L. C., Kazhdan, D., Singh, V., and Liò, P. GC-501 Explainer: Human-in-the-Loop Concept-based Explana-502 tions for Graph Neural Networks. In Workshop on Hu-503 man in the Loop Learning, volume 3. ICML, 2021. doi: 504 10.48550/ARXIV.2107.11889. URL https://arxiv. 505 org/abs/2107.11889.
- 507 Pope, P. E., Kolouri, S., Rostami, M., Martin, C. E., and 508 Hoffmann, H. Explainability methods for graph convolu-509 tional neural networks. In Proceedings of the IEEE/CVF 510 conference on computer vision and pattern recognition, 511 pp. 10772-10781, 2019. 512
- 513 Saha, S., Das, M., and Bandyopadhyay, S. GraphEx: A 514 User-Centric Model-Level Explainer for Graph Neural 515 Networks. In Maughan, K., Liu, R., and Burns, T. F. 516 (eds.), The First Tiny Papers Track at ICLR 2023, Tiny Pa-517 pers @ ICLR 2023, Kigali, Rwanda, May 5, 2023. Open-518 Review.net, 2023a. URL https://openreview. 519 net/pdf?id=CuE1F1M0_yR. 520
- 521 Saha, S., Das, M., and Bandyopadhyay, S. GraphEx: 522 A User-Centric Model-Level Explainer for Graph Neu-523 ral Networks. March 2023b. URL https:// 524 openreview.net/forum?id=CuE1F1M0 yR. 525
- 526 Sanfeliu, A. and Fu, K.-S. A distance measure between attributed relational graphs for pattern recognition. IEEE Transactions on Systems, Man, and Cybernetics, SMC-13 (3):353-362, 1983a. doi: 10.1109/TSMC.1983.6313167. 530
- 531 Sanfeliu, A. and Fu, K.-S. A distance measure between 532 attributed relational graphs for pattern recognition. IEEE 533 Transactions on Systems, Man, and Cybernetics, SMC-13 534 (3):353-362, May 1983b. ISSN 0018-9472, 2168-2909. 535 doi: 10.1109/TSMC.1983.6313167. URL http:// 536 ieeexplore.ieee.org/document/6313167/. 537
- 538 Schlichtkrull, M. S., De Cao, N., and Titov, I. Interpreting 539 graph neural networks for nlp with differentiable edge 540 masking. arXiv preprint arXiv:2010.00577, Proceedings 541 of International Conference on Learning Representations, 542 2020. 543
- 544 Schnake, T., Eberle, O., Lederer, J., Nakajima, S., Schütt, 545 K. T., Müller, K.-R., and Montavon, G. Higher-order 546 explanations of graph neural networks via relevant walks. 547 IEEE transactions on pattern analysis and machine intel-548 ligence, 44(11):7581-7596, 2021. 549

- Shin, Y.-M., Kim, S.-W., and Shin, W.-Y. PAGE: Prototype-Based Model-Level Explanations for Graph Neural Networks. 2022. doi: 10.48550/ARXIV.2210.17159. URL https://arxiv.org/abs/2210.17159.
- Veličković, P., Cucurull, G., Casanova, A., Romero, A., Liò, P., and Bengio, Y. Graph Attention Networks. 2017. doi: 10.48550/ARXIV.1710.10903. URL https://arxiv. org/abs/1710.10903.
- Vu, M. and Thai, M. T. Pgm-explainer: Probabilistic graphical model explanations for graph neural networks. Advances in neural information processing systems, 33: 12225-12235, 2020.
- Wang, X. and Shen, H.-W. GNNInterpreter: A Probabilistic Generative Model-Level Explanation for Graph Neural Networks. 2022. doi: 10.48550/ARXIV.2209.07924. URL https://arxiv.org/abs/2209.07924.
- Xu, K., Hu, W., Leskovec, J., and Jegelka, S. How Powerful are Graph Neural Networks? 2018. doi: 10.48550/ARXIV.1810.00826. URL https://arxiv. org/abs/1810.00826.
- Xuanyuan, H., Barbiero, P., Georgiev, D., Magister, L. C., and Liò, P. Global concept-based interpretability for graph neural networks via neuron analysis. In Proceedings of the Thirty-Seventh AAAI Conference on Artificial Intelligence and Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence and Thirteenth Symposium on Educational Advances in Artificial Intelligence, volume 37 of AAAI'23/IAAI'23/EAAI'23, pp. 10675-10683. AAAI Press, February 2023. ISBN 9781577358800. doi: 10.1609/aaai.v37i9.26267. URL https://doi.org/ 10.1609/aaai.v37i9.26267.
- Ying, Z., Bourgeois, D., You, J., Zitnik, M., and Leskovec, J. Gnnexplainer: Generating explanations for graph neural networks. Advances in neural information processing systems, 32, 2019.
- Yu, J., Xu, T., Rong, Y., Bian, Y., Huang, J., and He, R. Graph Information Bottleneck for Subgraph Recognition. October 2020. URL https://openreview.net/ forum?id=bM4Iqfq8M2k.
- Yuan, H., Tang, J., Hu, X., and Ji, S. XGNN: Towards Model-Level Explanations of Graph Neural Networks. In Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, pp. 430-438, Virtual Event CA USA, August 2020. ACM. ISBN 9781450379984. doi: 10.1145/3394486. 3403085. URL https://dl.acm.org/doi/10. 1145/3394486.3403085.

- Yuan, H., Yu, H., Wang, J., Li, K., and Ji, S. On explainability of graph neural networks via subgraph explorations. In *International conference on machine learning*, pp. 12241–
 12252. PMLR, 2021.
- Yuan, H., Yu, H., Gui, S., and Ji, S. Explainability in Graph
 Neural Networks: A Taxonomic Survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(5):
 5782–5799, May 2023. ISSN 1939-3539. doi: 10.1109/
 TPAMI.2022.3204236. URL https://ieeexplore.
 ieee.org/abstract/document/9875989.
- Zhang, Z., Liu, Q., Wang, H., Lu, C., and Lee, C. ProtGNN: Towards Self-Explaining Graph Neural Networks. Proceedings of the AAAI Conference on Artificial Intelligence, 36(8):9127–9135, June 2022. ISSN 2374-3468, 2159-5399. doi: 10.1609/aaai.v36i8. 20898. URL https://ojs.aaai.org/index. php/AAAI/article/view/20898.

605 A. Regularization Terms

606 Following the paradigm established by existing methods, an objective function for explanations typically contains two parts: 607 a term related to class prediction and a regularizer that enforces the generated explanations to be in-distribution. In XGNN, 608 the explanation generator is penalized during training for actions that violate manually-defined sets of rules, such as the 609 maximum number of bonds that can be formed with a certain atom in a molecule. In GNNInterpreter, the embedding of 610 the explanation graph needs to be close to the average embedding of graphs in the training set. While these regularization 611 strategies may help confine the explanation graph to a region of the input space where the model is well-defined, they 612 cannot guarantee the quality of the explanation. While regularization terms can normally be balanced through some 613 tuning procedure, this is impossible without knowing the ground-truth explanations for the GNN already, and attempting 614 to determine the weights by judging the generated graphs qualitatively increases the likelihood of mistakenly accepting 615 spurious explanations. Therefore, we do not apply any regularization in the objective function during our experiments, but 616 the proposed method is able to incorporate commonly used regularizers if desired. 617

B. Experimental Setup

618 619

638

639

620 In all experiments, the GNNs use GraphSAGE-style convolutions with sum being used as the aggregation operator, followed 621 by a global mean pooling layer, and finally several fully-connected (FC) layers. ReLU activations are placed between each 622 hidden layer. For the Is_Acyclic and Shapes datasets, the GNN uses 2 convolutional layers computing 16 features per node, 623 a FC layer computing 8 features, and a final FC layer to compute the class logits. For the MUTAG dataset, the GNN uses 2 624 convolutional layers computing 64 and 32 features per node, two FC layers computing 16 and 8 features per graph, and 625 a final FC layer to compute the logits. We implemented these GNNs using PyTorch-Geometric (Fey & Lenssen, 2019). 626 Models were trained for 200 epochs, optimizing with Adam (Kingma & Ba, 2015) with a learning rate of 10^{-3} and L2 627 regularization with weight 10^{-4} . 628

629 For the MUTAG dataset, XGNN's graph generator policy network was penalized when it violated valence constraints while 630 generating molecules, and no penalties were used on the other datasets. In the experiments with MIPExplainer, adjacency 631 matrices were constrained to be symmetric to represent undirected connected graphs without self-loops. For the MUTAG 632 dataset, node features were constrained to one-hot vectors by ensuring the sum of the elements in each row added up to 633 1. Any experiments lasting longer than 6 hours were automatically terminated, and we report the best solution found. To 634 ensure that any resemblance to target classes would not come from an initial solution, all runs for our method in 3 were 635 initialized with a path graph of n nodes. On the other hand, to test consistency, runs for MIPExplainer were initialized with 636 a graph generated by adding every possible edge to a line graph with probability 0.5. 637

C. Practical Considerations

In practice, it can be difficult to solve MIPs corresponding to large GNNs, and several techniques are needed to make the process tractable. Often, just finding an initial setting for all of the decision variables that satisfies all constraints is difficult. In our experiments, we found that this step can actually take longer than the subsequent optimization. This problem can be completely eliminated with a warm start. Starting from an arbitrary input graph (either from the dataset or not), we can compute a forward pass through the network to obtain a valid setting of initial values for almost all of the decision variables. In cases where additional constraints have been imposed on the graph, such as to ensure connectivity as described above, an input graph must be converted into the canonical form that also satisfies these constraints.

Floating-point precision errors can lead to serious problems for MIQP solvers, and in cases where decision variables can take both small and large values, a significant amount of time may be needed to avoid numerical instability. This is relevant when the weights of GNNs become very small, an effect often produced by regularization. However, we found that weights below a certain threshold (we chose 10^{-5}) could be floored to zero without significantly affecting the behavior of the network. All performance metrics for the networks used in the experiments were computed after the networks were pruned in this way. We also found that smoothing networks with regularization improved solution times.

654 Despite these measures, runtime remains the most significant drawback of the proposed approach. In our largest experiments, 655 we were not able to guarantee a global optimum within a single day. However, the largest reason for this runtime is the 656 amount of symmetry in our formulation. A single graph corresponds to a number of adjacency matrices that, in the worst 657 case, grows exponentially with the number of vertices it contains. This hinders our ability to tighten the upper bound while 658 exploring the search tree, since an existing global optimum may be transformed into another as we traverse a branch. In the 659 future, we plan to address this problem by introducing additional constraints to reduce the number of feasible adjacency
 matrices in each equivalence class defined by graph isomorphism. Despite the increased time needed to prove optimality,
 the proposed method often finds optimal solutions early in the search. Therefore, we impose a time limit during experiments
 to prove its practicality.

D. Extended Results



Figure 4. Improving solutions found while optimizing the explanation for the Wheel class from the Shapes dataset

			MIPExplainer	Average Edit Distan GNNInterpreter	ce XGNN
Dataset	Class	# Nodes			
Is_Acyclic	Acyclic	5	0.000 ± 0.000	3.000 ± 1.549	0.600 ± 0.490
		6	0.000 ± 0.000	3.200 ± 1.077	2.400 ± 1.200
		7	0.800 ± 0.980	3.000 ± 1.000	4.000 ± 1.414
		8	0.000 ± 0.000	4.600 ± 1.020	4.000 ± 0.775
	Cyclic	5	0.000 ± 0.000	2.600 ± 1.428	1.200 ± 0.600
		6	0.000 ± 0.000	2.200 ± 1.249	2.800 ± 1.077
		7	0.000 ± 0.000	3.000 ± 1.000	3.400 ± 1.744
		8	0.000 ± 0.000	4.600 ± 1.020	5.600 ± 1.497
MUTAG	Mutagen	5	0.000 ± 0.000	7.000 ± 1.483	5.200 ± 1.327
		6	0.000 ± 0.000	7.300 ± 1.616	7.000 ± 1.414
		7	2.800 ± 3.429	8.500 ± 1.285	8.500 ± 1.500
		8	0.000 ± 0.000	9.600 ± 2.458	10.300 ± 1.345
	Nonmutagen	5	0.000 ± 0.000	5.800 ± 1.327	4.600 ± 1.428
		6	0.000 ± 0.000	7.500 ± 1.204	7.000 ± 1.549
		7	1.600 ± 1.960	8.400 ± 1.428	7.700 ± 0.781
		8	0.000 ± 0.000	9.500 ± 1.025	10.400 ± 2.107
Shapes	Grid	5	0.000 ± 0.000	4.400 ± 2.107	1.200 ± 0.748
		6	0.000 ± 0.000	2.600 ± 0.917	1.000 ± 0.632
		7	0.000 ± 0.000	3.000 ± 0.894	3.400 ± 1.855
		8	0.000 ± 0.000	4.200 ± 0.872	6.200 ± 2.441
	Lollipop	5	0.000 ± 0.000	2.800 ± 1.249	1.600 ± 1.428
		6	0.800 ± 0.980	3.000 ± 0.894	2.000 ± 0.894
		7	0.000 ± 0.000	4.000 ± 1.000	4.000 ± 1.000
		8	0.000 ± 0.000	4.400 ± 1.356	6.000 ± 2.145
	Star	5	0.000 ± 0.000	5.200 ± 2.272	1.200 ± 0.980
		6	0.000 ± 0.000	2.600 ± 0.663	1.000 ± 0.632
		7	0.000 ± 0.000	3.200 ± 0.748	3.400 ± 1.200
		8	0.000 ± 0.000	4.600 ± 1.356	6.600 ± 3.611
	Wheel	5	0.000 ± 0.000	1.800 ± 0.980	0.600 ± 0.490
		6	0.000 ± 0.000	4.000 ± 1.183	2.000 ± 0.775
		7	1.200 ± 0.600	3.400 ± 1.428	5.200 ± 2.750
		8	1.800 ± 0.600	4.000 ± 0.894	8.000 ± 3.924

Table 5. Average edit distance between 5 generated example graphs. The time limit was one half hour for all experiments. MIPExplainer had the lowest average edit distance in all experiments.



Figure 5. Solver metrics for 5 runs explaining the Mutagen class of MUTAG (Top) and the Acyclic class of Is_Acyclic (bottom) with 5 nodes: On the left, the current best solution's objective (blue) and upper bound (red) converging to the same global optimum (the dotted black line). On the right, the number of explored (green)/unexplored (orange) nodes during the search.

Dataset	Class	# Nodes	MIPExplainer	Runtime GNNInterpreter	XGNN
Is_Acyclic	Acyclic	5	3.965 ± 0.313	7.610 ± 0.015	10.040 ± 0.13
•		6	13.752 ± 0.767	7.626 ± 0.022	13.482 ± 0.52
		7	119.144 ± 19.283	7.679 ± 0.054	15.419 ± 0.50
		8	919.366 ± 387.422	7.684 ± 0.013	17.265 ± 0.4
	Cyclic	5	1.701 ± 0.356	0.027 ± 0.008	9.842 ± 0.05
	-	6	1.781 ± 0.521	0.024 ± 0.000	12.877 ± 0.4
		7	1.902 ± 0.102	0.030 ± 0.003	15.257 ± 0.3
		8	1.945 ± 0.098	0.115 ± 0.051	18.216 ± 1.2
MUTAG	Mutagen	5	312.586 ± 120.643	0.024 ± 0.001	8.701 ± 0.45
		6	901.010 ± 82.795	0.025 ± 0.004	10.305 ± 0.2
		7	1157.387 ± 386.822	0.060 ± 0.060	12.069 ± 0.2
		8	3031.742 ± 288.508	0.114 ± 0.007	14.727 ± 0.6
	Nonmutagen	5	1916.024 ± 753.053	5.633 ± 5.109	8.747 ± 0.53
		6	7201.688 ± 0.890	7.599 ± 4.213	10.403 ± 0.3
		7	7201.532 ± 0.087	9.659 ± 0.049	12.595 ± 0.7
		8	7201.842 ± 0.122	7.857 ± 4.304	13.946 ± 0.4
Shapes	Grid	5	4.964 ± 1.060	7.563 ± 0.018	11.644 ± 0.1
		6	17.508 ± 3.246	7.614 ± 0.040	14.964 ± 0.1
		7	92.649 ± 19.726	7.648 ± 0.022	17.817 ± 0.1
		8	1886.365 ± 537.301	7.667 ± 0.007	20.497 ± 0.7
	Lollipop	5	7.560 ± 1.296	7.631 ± 0.016	11.543 ± 0.1
		6	60.700 ± 4.313	7.634 ± 0.013	15.048 ± 0.4
		7	172.511 ± 44.716	7.661 ± 0.019	17.934 ± 0.3
		8	4261.094 ± 2267.636	7.705 ± 0.027	19.852 ± 0.3
	Star	5	4.050 ± 0.711	7.560 ± 0.014	11.512 ± 0.1
		6	18.200 ± 2.062	7.592 ± 0.007	14.955 ± 0.2
		7	176.005 ± 113.007	7.624 ± 0.014	17.942 ± 0.4
		8	260.384 ± 40.326	7.659 ± 0.016	19.931 ± 0.6
	Wheel	5	4.256 ± 1.225	7.716 ± 0.255	11.550 ± 0.0
		6	107.301 ± 134.613	7.641 ± 0.010	15.507 ± 0.7
		7	103.210 ± 8.210	7.652 ± 0.010	18.047 ± 0.3
		8	5576.961 ± 1071.344	7.706 ± 0.051	20.095 ± 0.3

Table 6. Runtime of explanation methods given a time limit of 2 hours, averaged over 5 runs

Submission and Formatting Instructions for ICML 2024

			Is_Acyclic Output Logits				
			MIPExplainer	GNNInterpreter	XGNN		
Class	# Nodes		-	-			
Acyclic	5	Cyclic Logit	-8.981 ± 0.000	-1.343 ± 4.048	6.417 ± 0.7		
		Acyclic Logit	11.896 ± 0.000	1.835 ± 5.022	-7.456 ± 0.8		
	6	Cyclic Logit	-9.722 ± 0.000	1.452 ± 4.292	7.333 ± 1.8		
		Acyclic Logit	12.932 ± 0.000	$\textbf{-1.577} \pm 5.222$	-8.490 ± 2.1		
	7	Cyclic Logit	-9.184 ± 2.476	4.607 ± 1.646	6.520 ± 1.7		
		Acyclic Logit	12.247 ± 3.341	$\textbf{-5.415} \pm \textbf{1.857}$	-7.573 ± 1.9		
	8	Cyclic Logit	-10.754 ± 0.000	4.924 ± 1.633	6.372 ± 1.8		
		Acyclic Logit	14.409 ± 0.000	$\textbf{-5.771} \pm \textbf{1.843}$	$-7.400 \pm 2.$		
Cyclic	5	Cyclic Logit	7.194 ± 0.000	-5.769 ± 2.112	5.665 ± 1.0		
		Acyclic Logit	$\textbf{-8.333} \pm 0.000$	7.516 ± 2.844	$-6.609 \pm 1.$		
	6	Cyclic Logit	10.260 ± 0.000	0.200 ± 4.304	6.997 ± 2.1		
		Acyclic Logit	-11.793 ± 0.000	$\textbf{-0.080} \pm 5.363$	-8.111 ± 2.12		
	7	Cyclic Logit	13.488 ± 0.000	-0.293 ± 1.360	7.080 ± 1.0		
		Acyclic Logit	-15.436 ± 0.000	0.524 ± 1.625	-8.205 ± 1.5		
	8	Cyclic Logit	16.870 ± 0.000	4.585 ± 1.148	7.587 ± 1.1		
		Acyclic Logit	-19.184 ± 0.000	-5.328 ± 1.362	-8.777 ± 1.2		

Table 7. Logits for Is_Acyclic explanation graphs generated by MIPExplainer, averaged over 5 runs with random initial solutions

			Ν	ts	
			MIPExplainer	GNNInterpreter	XGNN
Class	# Nodes				
Mutagen	5	Nonmutagen Output Logit	-19.249 ± 0.000	-0.422 ± 5.318	3.825 ± 5.060
		Mutagen Output Logit	11.269 ± 0.000	4.048 ± 7.581	-0.041 ± 5.50
	6	Nonmutagen Output Logit	-30.439 ± 0.000	-7.649 ± 7.386	-1.738 ± 8.14
		Mutagen Output Logit	17.673 ± 0.000	18.074 ± 13.474	7.321 ± 8.283
	7	Nonmutagen Output Logit	-43.654 ± 0.000	-9.854 ± 13.851	-1.755 ± 3.12
		Mutagen Output Logit	25.229 ± 0.000	20.275 ± 23.446	4.828 ± 6.00
	8	Nonmutagen Output Logit	-59.041 ± 0.000	-21.218 ± 12.229	-1.859 ± 5.72
		Mutagen Output Logit	34.029 ± 0.000	41.498 ± 19.722	5.446 ± 7.36
Nonmutagen	5	Nonmutagen Output Logit	6.691 ± 0.000	1.887 ± 4.245	-3.481 ± 1.84
		Mutagen Output Logit	$\textbf{-7.612}\pm0.000$	-0.107 ± 5.117	8.834 ± 2.69
	6	Nonmutagen Output Logit	6.830 ± 0.000	7.146 ± 8.787	2.359 ± 3.89
		Mutagen Output Logit	$\textbf{-7.779}\pm0.000$	-0.851 ± 6.609	1.204 ± 6.98
	7	Nonmutagen Output Logit	7.185 ± 0.098	-3.194 ± 7.405	-0.841 ± 6.12
		Mutagen Output Logit	$\textbf{-8.174} \pm 0.102$	10.138 ± 9.418	3.811 ± 8.57
	8	Nonmutagen Output Logit	7.236 ± 0.000	-5.436 ± 10.869	-9.350 ± 15.2
		Mutagen Output Logit	-8.251 ± 0.000	13.410 ± 16.078	19.037 ± 22.0

Table 8. Logits for MUTAG explanation graphs generated by MIPExplainer, averaged over 5 runs with random initial
 solutions

887						
888				S	hapes Output Logits	(1)
889				MIPExplainer	GNNInterpreter	XGNN
890	Class	# Nodes				
891	Grid	5	Lollipop Logit	-1.343 ± 0.000	-10.019 ± 6.822	-5.121 ± 1.118
892		-	Wheel Logit	-19.762 ± 0.000	-38.818 ± 40.875	3.259 ± 6.784
893			Grid Logit	9.577 ± 0.000	-8.103 ± 23.501	8.262 ± 0.688
894			Star Logit	-20.261 ± 0.000	-14.756 ± 18.554	-31.471 ± 4.109
895		6	Lollipop Logit	-0.846 ± 0.000	-3.386 ± 3.512	-4.907 ± 1.180
896			Wheel Logit	-1.949 ± 0.000	-21.545 ± 11.607	5.398 ± 3.124
897			Grid Logit	9.154 ± 0.000	-6.520 ± 11.184	4.816 ± 2.687
898			Star Logit	-31.297 ± 0.000	-17.602 ± 8.887	-30.366 ± 3.257
899		7	Lollipop Logit	-1.269 ± 0.000	-2.233 ± 2.319	-1.976 ± 0.763
900			Wheel Logit	-11.014 ± 0.000	-5.299 ± 5.715	-4.893 ± 2.458
901			Grid Logit	9.972 ± 0.000	3.110 ± 3.284	-1.793 ± 7.365
902			Star Logit	-25.884 ± 0.000	-23.110 ± 2.283	-21.878 ± 3.027
903		8	Lollipop Logit	-1.502 ± 0.000	-2.884 ± 3.272	-2.444 ± 1.782
904			Wheel Logit	$\textbf{-8.745} \pm 0.000$	-11.408 ± 6.947	-5.672 ± 8.371
905			Grid Logit	9.270 ± 0.000	1.346 ± 5.852	-7.951 ± 13.444
900			Star Logit	-26.275 ± 0.000	-19.579 ± 3.193	-19.534 ± 7.940
907	Lollipop	5	Lollipop Logit	-2.791 ± 0.000	-8.329 ± 8.233	-5.837 ± 0.697
908			Wheel Logit	-25.569 ± 0.000	-47.927 ± 27.593	-0.149 ± 22.146
909			Grid Logit	$\textbf{-5.128} \pm 0.000$	-15.098 ± 19.693	3.951 ± 10.263
910			Star Logit	-19.403 ± 0.000	-9.227 ± 13.375	-30.313 ± 12.808
911		6	Lollipop Logit	1.103 ± 1.065	-3.452 ± 3.391	-2.709 ± 1.291
912			Wheel Logit	-17.920 ± 8.127	-22.403 ± 11.090	-1.503 ± 3.706
915			Grid Logit	-7.984 ± 8.479	-3.795 ± 8.878	2.021 ± 0.575
914			Star Logit	-18.629 ± 5.279	-15.728 ± 7.512	-25.136 ± 2.219
915		7	Lollipop Logit	2.700 ± 0.000	-3.375 ± 2.803	-3.732 ± 0.426
017			Wheel Logit	-18.050 ± 0.000	-11.131 ± 14.287	-1.052 ± 5.294
917			Grid Logit	$\textbf{-5.612} \pm 0.000$	$\textbf{-4.683} \pm \textbf{8.943}$	1.356 ± 2.642
010			Star Logit	-21.298 ± 0.000	-22.799 ± 8.634	-24.479 ± 4.387
920		8	Lollipop Logit	7.224 ± 0.000	-2.637 ± 1.622	-2.184 ± 3.039
021			Wheel Logit	-15.080 ± 0.000	-19.483 ± 17.827	-13.259 ± 5.576
077			Grid Logit	$\textbf{-3.129}\pm0.000$	1.137 ± 6.363	-6.440 ± 5.380
923			Star Logit	-19.795 ± 0.000	-18.262 ± 9.027	-14.606 ± 4.264
140						

Table 9. Logits for explanation graphs generated by MIPExplainer for the "Grid" and "Lollipop" classes of the shapes dataset, averaged over 5 runs with random initial solutions

942		
943 Sh	apes Output Logits	(2)
944 MIPExplainer	GNNInterpreter	XGNN
945 Class # Nodes		
946 Star 5 Lollipop Logit -16.772 ± 0.000	-12.575 ± 12.292	-5.020 ± 1.008
947 Wheel Logit -67.715 ± 0.000	-42.824 ± 40.005	2.464 ± 5.862
Grid Logit -72.108 ± 0.000	-31.701 ± 46.376	7.990 ± 0.654
Star Logit 12.894 ± 0.000	-8.555 ± 20.510	-30.609 ± 2.926
6 Lollipop Logit -13.064 ± 0.000	-3.980 ± 2.651	-4.822 ± 1.125
951 Wheel Logit -64.317 ± 0.000	-21.691 ± 9.404	5.352 ± 3.101
952 Grid Logit -48.784 ± 0.000	-3.585 ± 12.975	4.037 ± 2.655
953 Star Logit 12.183 ± 0.000	-15.661 ± 2.185	-29.810 ± 3.144
954 7 Lollipop Logit -11.148 ± 0.000	-2.507 ± 2.511	-1.469 ± 2.153
Wheel Logit -61.412 ± 0.000	-22.296 ± 21.034	-9.213 ± 5.345
Grid Logit -33.850 ± 0.000	-3.210 ± 13.723	-6.549 ± 8.455
957 Star Logit 13.072 ± 0.000	-18.741 ± 8.891	-18.000 ± 4.776
958 8 Lollipop Logit -10.722 ± 0.000	-3.702 ± 1.427	-1.924 ± 1.894
Wheel Logit -58.630 ± 0.000	-7.587 ± 15.924	$\textbf{-8.287} \pm \textbf{9.256}$
Grid Logit -23.479 ± 0.000	$\textbf{-1.624} \pm \textbf{8.923}$	-11.054 ± 12.549
Star Logit 14.536 ± 0.000	-21.715 ± 8.531	-18.099 ± 8.888
Wheel 5 Lollipop Logit -6.180 ± 0.000	-7.290 ± 7.960	-5.957 ± 0.275
Wheel Logit 13.062 ± 0.000	-42.230 ± 28.326	8.335 ± 2.176
Grid Logit 10.323 ± 0.000	-3.870 ± 18.922	8.057 ± 0.745
Star Logit -40.319 ± 0.000	-12.452 ± 10.727	-34.470 ± 2.362
6 Lollipop Logit -5.616 ± 0.000	-4.918 ± 4.503	-3.865 ± 1.889
Wheel Logit 6.921 ± 0.000	-18.368 ± 31.653	1.860 ± 6.420
Grid Logit 2.541 ± 0.000	3.615 ± 6.467	3.038 ± 3.642
Star Logit -29.458 ± 0.000	-22.223 ± 14.202	-27.369 ± 5.667
7 Lollipop Logit -5.328 ± 0.402	-4.489 ± 2.121	-1.477 ± 4.376
Wheel Logit 7.040 ± 1.338	-5.192 ± 15.185	-8.111 ± 10.727
Grid Logit 2.960 ± 1.279	5.168 ± 3.163	-4.639 ± 12.694
Star Logit -29.519 ± 1.233	-24.700 ± 9.701	-20.903 ± 8.152
8 Lollipop Logit -4.147 ± 1.128	-2.193 ± 0.699	-1.513 ± 2.686
Wheel Logit 4.468 ± 3.188	$\textbf{-5.022} \pm \textbf{6.210}$	-11.867 ± 4.893
Grid Logit 1.505 ± 0.752	-5.701 ± 13.464	-3.223 ± 12.412
Star Logit -27.293 ± 1.755	-20.266 ± 6.735	-17.747 ± 6.596

Table 10. Logits for explanation graphs generated by MIPExplainer for the "Star" and "Wheel" classes of the shapes dataset, averaged over 5 runs with random initial solutions