Explaining Graph Neural Networks with Mixed-Integer Programming

Anonymous CVPR submission

Paper ID 16514

Abstract

001 Graph Neural Networks (GNNs) have become a popu-002 lar tool for learning from graph-structured data. While 003 they provide state-of-the-art graph learning performance, their lack of transparency hinders our ability to trust and 004 005 understand them, ultimately limiting the areas where they can be applied. Many explanation methods exist to ex-006 007 plain individual predictions made by GNNs, but there are 800 fewer ways to gain more general insight into the patterns they have been trained to identify. Most proposed meth-009 010 ods for model-level GNN explanations attempt to generate graphs that exemplify these patterns, but the discreteness 011 012 of graphs and the nonlinearity of deep GNNs make finding 013 such graphs difficult. In this paper, we formulate the search 014 for an explanatory graph as a mixed-integer programming (MIP) problem, in which decision variables specify the ex-015 planation graph and the objective function represents the 016 017 quality of the graph as an explanation for a GNN's predic-018 tions of an entire class in the dataset. This approach, which 019 we call MIPExplainer, allows us to directly optimize over 020 the discrete input space and find globally optimal solutions with a minimal number of hyperparameters. MIPExplainer 021 022 is able to find accurate explanations on both synthetic and 023 real-world datasets, and outperforms existing methods for 024 generating explanations for GNNs.

025 **1. Introduction**

Graph neural networks (GNNs), such as graph convolu-026 027 tional networks (GCN) [12], GraphSAGE [8], and graph attention networks (GAT) [26], provide a family of pow-028 029 erful tools for modelling graphs that learn from both the features contained in nodes and edges and the structure of 030 031 the graph itself. While they have achieved state-of-the-art performance across a wide range of application domains, 032 their applications are still limited by our inability to explain 033 their behavior. Without being able to explain the patterns 034 GNNs rely on to make predictions, it is impossible to jus-035 036 tify their use in applications where trust and safety are im-037 portant, and there is no way to extract useful information from them. This has motivated a significant amount of research into techniques for GNN explainability. 039

Although there is research into inherently explainable 040 deep learning, such as ProtGNN [35] and GIB [31] for 041 GNNs, this approach generally comes at the cost of perfor-042 mance. Therefore, most explanation methods are designed 043 to be applied post-hoc to trained models. Elucidating how 044 a GNN "reasons" on graphs presents a completely new set 045 of challenges from deep learning explainability over image 046 and text data. Encodings of graphs lack the same spatial 047 information, and their structure must be represented explic-048 itly, for example by adjacency matrices or edge lists. These 049 representations are inherently discrete, which violates as-050 sumptions necessary for the application of commonly-used 051 gradient based methods [5, 19]. The outputs of GNNs are 052 further expected to be permutation equivariant with respect 053 to the ordering of the nodes, which further imposes con-054 straints on generating explanation graphs for a GNN model. 055 A one-size-fits-all strategy may not work with the wide va-056 riety of existing graph types (e.g., directed and undirected 057 graphs, graphs with node and edge features) and GNN lay-058 ers (e.g., convolutions in GCNs, GATs, GINs). Although 059 it can be further extended, our method will focus on un-060 weighted graphs and GraphSAGE-style graph convolutions. 061

1.1. Related Work

GNN interpretation has been largely focused on instance-063 level explanation which aims to explain the reasoning be-064 hind individual predictions. As identified in [34], at least six 065 categories of instance-level GNN explanation methods have 066 been proposed so far: gradient-based [20], perturbation-067 based [18, 23, 30, 33], surrogate [27], generation-based 068 [13], decomposition [24], and counterfactual-based [17] 069 methods. These methods do not immediately provide in-070 sights into the overall patterns a GNN has identified to make 071 predictions, but it is possible to consolidate instance-level 072 explanations to reveal model-level patterns. For example, 073 we can employ purely statistical methods (e.g., a rank-sum 074 test) to determine whether there are nodes/edges shared by 075 076 a significant portion of the explanations. A more recent technique is GLGExplainer [2], which finds smaller com-077

174

078 ponents of the extracted explanations that can be used to build logical expressions consistent with the GNN's predic-079 080 tions (e.g. the GNN predicts a class of graphs given the presence of a certain subgraph and the absence of another 081 082 subgraph). However, these methods are limited by the scope of the training data, and can be influenced by bias in the 083 dataset. A direct model-level explanation can offer more 084 085 faithful explanations, and is more useful for determining the 086 degree of bias in the model itself. In this work, we focus on 087 post-hoc model-level explainability for GNNs.

Relatively few methods exist to explain GNNs at the 088 089 model level. XGNN [32] is the most widely used, and 090 serves as the only baseline in several recent papers that focus on similar objectives [2, 21, 25, 28]. These methods aim 091 092 to generate graphs that exemplify graph structures used by a trained GNN for making classifications, without straying 093 094 too far from the distribution of the training data where the 095 model is not well-defined. XGNN involves training a second neural network by reinforcement learning to generate 096 097 graphs that maximize the original GNN's prediction for a specific class. OrphicX [15] and GEM [14] are two other 098 099 methods that involve training a second neural network to 100 generate explanations, but they use different objective func-101 tions designed to encourage causal explanations. The use of 102 a second black box model in these methods reduces their reliability and makes interpreting the generated explanations 103 significantly more difficult. 104

105 Other recent methods avoid training a second neural network in various ways. GNNInterpreter [28] and GraphEx 106 107 [22] assume that the graphs in the dataset are sampled from a set of underlying distributions parameterized by continu-108 ous latent parameters. GNNInterpreter defines an objective 109 function similar to XGNN during training, maximizing a 110 111 target class's logit while penalizing the distance between the 112 embedding of the generated graph and the mean embedding 113 of the training data to keep explanations in-distribution, and learns parameters through Monte Carlo gradient estima-114 115 tion. GraphEx attempts to solve this problem using max-116 imum likelihood estimation, which does not require gradients from the model and only relies on its predictions of 117 graphs in the dataset. PAGE [25] also maximizes predicted 118 class probability, but restricts explanations to common sub-119 graphs in the target class. 120

In the past, discrete optimization techniques have been 121 122 applied to deep neural networks to solve inverse problems, such as in [3] and [1]. Finding a graph that maximizes the 123 124 GNN output is a type of inverse problem. However, these 125 existing methods were only defined for layers consisting of a linear transformation with ReLU activations, and rely on 126 costly automated bound-tightening procedures to achieve 127 feasible runtimes. This work will reformulate the optimiza-128 129 tion problem for explainability, generalize the application of 130 mixed-integer programming on standard neural networks to a range of GNN architectures, and derive explicit bounds131for continuous variables that can be calculated with a single132forward pass through the network.133

1.2. Current Issues and Our Contributions

There are several common problems among existing ap-135 proaches for model-level GNN explanation. For one, they 136 all have many hyperparameters that can change the gener-137 ated explanation graphs. Generating a high-quality expla-138 nation requires specific settings for these hyperparameters, 139 but the more times a user runs the same algorithm with dif-140 ferent hyperparameters, the more likely they are to gener-141 ate spurious explanations. In cases where hyperparameters 142 modify the objective function used to calculate the quality 143 of explanations (e.g. weights for regularization terms), it is 144 impossible to quantitatively compare the quality of the gen-145 erated explanations, even though they may be extremely dif-146 ferent. This problem is exasperated by the fact that several 147 of these methods rely on stochastic gradient optimization to 148 optimize their explanations. Due to the stochastic nature of 149 this approach and the necessity of setting a maximum num-150 ber of iterations, the objective value might be far from the 151 global optimum when an explanation graph is generated, 152 without any way of knowing that this is the case. Thus, the 153 corresponding explanation may not represent the informa-154 tion that the objective function was designed to extract. 155

In this paper, we propose a new explanation method 156 based on mixed-integer programming (MIP) for finding 157 class-representative input graphs. This has several benefits 158 over existing approaches. (1) Unlike any existing method, 159 we can directly optimize over the discrete space of possi-160 ble input graphs, without any restrictions on types of node 161 and edge features. The only assumption we make about 162 the space of graphs is a bounded size, and we do not re-163 quire any assumptions about its underlying distribution. (2) 164 This approach has a minimal number of hyperparameters 165 that influence the explanation (only the size of the expla-166 nation must be specified), facilitating the application of our 167 approach and mitigating the effects of bias when analyzing 168 the results. (3) We prove that our MIP formulation has a 169 globally optimal solution, and in many cases we can find 170 and verify this solution. In cases where this is intractable, 171 we can place an upper bound on the optimal solution, guar-172 anteeing the quality of the generated explanation. 173

2. MIPExplainer

Our model-level explanation also seeks to optimize an input graph G = (X, A) so that the GNN predicts a certain class with the highest possible probability, where X contains the d attributes for each node as row vectors for N nodes in the graph and $A = (a_{ij})$ represents the N by N adjacency matrix. We focus on the case of a binary adjacency matrix where $a_{ij} \in \{0, 1\}$, so that $a_{ij} = 1$ indicates there is an 175 176 177 178 178 179 180 182 edge between nodes *i* and *j*. Let a GNN realize a function 183 $f_c(G, \theta)$ that maps *G* to the probabilities of several classes 184 indexed by *c* and θ contains all the learned parameters in the 185 GNN. During the GNN training, *G* is given and θ needs to 186 be determined, whereas in many model explanation meth-187 ods, θ has been fixed, and we optimize *G* (i.e., *X* and *A*) to 188 maximize $f_c(G, \theta)$ (or a related objective).

The proposed MIP will optimize G in terms of the val-189 190 ues of A and X. Each layer of the GNN imposes a set of constraints in the MIP. We add decision variables to repre-191 sent the output of each layer, which are constrained to the 192 correct values in terms of the decision variables represent-193 ing the output of the previous layer. For example, for a fully 194 connected layer, a new matrix of decision variables Y' will 195 be added to the model and constrained with Y' = WY' + b, 196 where Y' are the decision variables representing the outputs 197 of a previous layer and W and b are the model's learned pa-198 199 rameters. Since the outputs of subsequent layers are con-200 strained exactly, ultimately all of the constraints define the 201 feasible region of X and A. We place constraints on nodes 202 and edges (or entries of A) so that the derived explanation forms a connected graph with valid features, and we can 203 further constrain X and A to reduce the number of candi-204 205 date solutions for a single graph since the GNN is permu-206 tation equivariant with respect to the order of nodes. In the 207 subsequent sections, we provide a detailed derivation of our MIP formulation by discussing the objective function and 208 the various constraints. 209

210 2.1. Objective Function

Following the paradigm established by existing methods, an 211 212 objective function for explanations typically contains two 213 parts: a term related to class prediction and a regularizer that enforces the generated explanations to be in-distribution. 214 215 In XGNN, the explanation generator is penalized during 216 training for actions that violate manually-defined sets of 217 rules, such as the maximum number of bonds that can be 218 formed with a certain atom in a molecule. In GNNInterpreter, the embedding of the explanation graph needs to be 219 close to the average embedding of graphs in the training set. 220 221 While these regularization strategies may help confine the explanation graph to a region of the input space where the 222 223 model is well-defined, they cannot guarantee the quality of the explanation. While regularization terms can normally 224 225 be balanced through some tuning procedure, this is impossible without knowing the ground-truth explanations for 226 227 the GNN already, and attempting to determine the weights 228 by judging the generated graphs qualitatively increases the likelihood of mistakenly accepting spurious explanations. 229 Therefore, we do not apply any regularization in the ob-230 jective function during our experiments, but the proposed 231 232 method is able to incorporate commonly used regularizers 233 if desired.



Figure 1. Logits for the Star and Wheel Graphs in the Shapes Dataset

While maximizing a single logit while disregarding the 234 logits of other classes in the denominator (e.g., as done by 235 GNNInterpreter) is possible, this may lead to low quality 236 explanations in some circumstances. Predictions are made 237 based on the difference between the logits, and the absolute 238 value of a single logit may be unrelated to the prediction 239 of the network. To illustrate this, after training a GNN to 240 classify star graphs and wheel graphs of varying sizes (a 241 task defined in [28]), we plotted the logits it assigned to 242 the training data for both classes in Figure 1. Note that the 243 maximum logit for the wheel class is actually assigned to 244 a correctly-classified star graph. Thus, simply maximizing 245 the logit for wheels will not produce an effective explana-246 tion for the wheel class. 247

To accurately find class-discriminative information, we 248 should maximize the difference between the logit of the tar-249 get class and the logit of the other classes. Maximizing the 250 normalized probability, as done by XGNN, is possible but 251 can lead to numerical instability when maximizing because 252 253 improvements get exponentially smaller as the magnitude of the logits increases. We can form an objective function as 254 a linear combination of all logits but with a positive coeffi-255 cient only for the target class, but at the risk of a single neg-256 ative logit dominating the objective value. To mitigate this 257 problem, we can maximize the difference between the logit 258 of the target class and the maximum of the other classes. In 259 our observation, the latter is more effective, so we focus on 260 discussing the following objective function: 261

$$\max_{G} \left(f_c(G,\theta) - \max_{i \neq c}(f_i(G,\theta)) \right), \quad (1) \quad 262$$

263

264

265

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.266
267We 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.268

309

314

272If GraphSAGE convolution layers are employed in the273GNN, the updated node representations X' after each layer274are calculated from existing node representations X with275the formula

276
$$X' = \sigma(XW_1 + \text{Aggregation}(X)W_2 + b)$$
(2)

where the aggregation is commonly realized by summing over neighbors' attributes, i.e., Aggregation(X) = AX.

279 Assume that a GNN model has L_c GraphSAGE-based convolution layers with sum aggregations and ReLU activa-280 tions, followed by a global feature-wise sum pooling layer 281 and L_f fully connected (FC) layers with ReLU activations. 282 In total, there are $L = L_c + 1 + L_f$ layers (indexed as 283 $L_i, i \in 1, ..., L$). We will use the following notation: $W_j^{(i)}$, matrices of scalars, denote the GNN's matrices of learned 284 285 weights in layer *i*. The vectors of scalars $b^{(i)}$ denote the 286 GNN's learned bias vectors in layer *i*. For notation con-287 venience, we also denote $X^{(0)} = X$, where x_{ij} is the *j*th 288 feature of node *i*. 289

We will also add the following decision variables to our 290 formulation, which will be constrained to the correct values 291 based on X and A: $\Phi^{(i)}$ represents the output of layer *i* be-292 fore the activation function, $X^{(i)}$ represents $ReLU(\Phi^{(i)})$, 293 the output of layer *i*. We also represent $ReLU(-\Phi^{(i)})$ by 294 $B^{(i)}$, while $Z^{(i)}$ are binary indicators representing the truth 295 value of $\Phi^{(i)} > 0$ elementwise¹. d is an indicator vector 296 representing whether each element of $\Phi^{(L)}$ is the maximum 297 element in the output of layer L not including the target 298 class (i.e., dimension j of d is 1 when dimension j of $\Phi^{(L)}$ 299 is the maximum, while the rest are all 0's), and y represents 300 the value of the maximum output of the GNN that is not for 301 the target class. Aside from the binary variables A, $Z^{(i)}$, 302 and d, all other variables are continuous. 303

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

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

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

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

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

311
$$\Phi^{(i)} = X^{(i-1)} W_1^{(i)} + b^{(i)}.$$
 (5)

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

$$X^{(i)} - B^{(i)} = \Phi^{(i)}, \tag{6} 315$$

$$X^{(i)} \le MZ^{(i)},$$
 (7) 316

$$B^{(i)} \le M(1 - Z^{(i)}),$$
 (8) 317

$$0 \le X^{(i)}, B^{(i)} \le M. \tag{9} 318$$

For the pooling layer, we simply have that 319 $X^{(L_c+1)} = \Phi^{(L_c+1)}$. To constrain d_j and y: 320

$$y \ge X_{\neq c}^{(L)},\tag{10} \qquad \textbf{321}$$

$$y \le X_{\neq c}^{(L)} + (\max(U_{X_{\neq c}^{(L)}})\mathbf{1} - L_{X_{\neq c}^{(L)}})(1-d), \quad (11) \quad 322$$

$$\sum_{j} d_j = 1, d_j \in \{0, 1\},$$
(12) 323

where $L_{X_{\neq a}^{(L)}}$ and $U_{X_{\neq a}^{(L)}}$ represent lower and upper bounds 324 for the decision variables in $X^{(L)}$ excluding the output of 325 class c. A method to calculate these bounds based on the 326 bounds of the input will be discussed in a later section. 327 Most of these constraints are linear in terms of the deci-328 sion variables except Eq.(3) where decision variables A and 329 $X^{(i)}$ multiply to form quadratic terms. Because A is bi-330 nary, these terms can be equivalently reformulated into lin-331 ear functions. There are several ways to perform the lin-332 earization of quadratic terms with both continuous and bi-333 nary variables. We describe one such method by change of 334 variables [10]. For a given binary variable $a \in A$ and a con-335 tinuous variable $x \in X^{(i)}$ bounded by M, let $e = a \times x$ be a 336 new intermediate decision variable. Let $E^{(i)}$ be the matrix 337 of $AX^{(i)}$ where entries are all calculated by summing the 338 corresponding e's. Constraints in Eq.(3) can be rewritten as 339 follows with additional bound constraints: 340

$$\Phi^{(i)} = X^{(i-1)} W_1^{(i)} + E^{(i)} W_2^{(i)} + b^{(i)}, \qquad (13) \qquad \textbf{341}$$

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

350

$$x - M(1 - a) \le e \le x + M(1 - a).$$
 (15) 343

Now, our MIP has been transformed into a problem that344maximizes Eq.(1) which can be calculated as $X_c^{(L)} - y$ subject to constraints Eq.(13-15) and Eq.(4-12). Note that this345MIP has a concave objective function and all linear constraints when integrality is relaxed, so it is a mixed-integer348linear program (MILP).349

2.3. Constraints on *A* and *X*

Additional constraints can be placed on A and X when gen-
erating 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 the351
352

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

392

393

394

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 360 $node_1, \ldots, node_n$ to ensure that the explanation graph is 361 connected, or in the case of directed graphs, weakly con-362 363 nected (i.e. connected ignoring the directionality of the 364 edges). We require that there is at least one edge between node_i and the set of nodes {node_i | i > j}. This be-365 comes a constraint on A, and specifically for each i with 366 $0 \le i < n$, we add the constraint $\sum_{\{j \mid i > j\}} (a_{ij} + a_{ji}) \ge 1$. 367 These constraints also partially alleviate the effect of equiv-368 369 ariance (where the same graph can have many different A), because they reduce the MIP's feasible region, but not the 370 set of candidate graphs. This can easily be proven by show-371 ing that for any (weakly) connected graph, the nodes can be 372 ordered in a way that satisfies these constraints by running 373 374 depth-first-search (DFS) on such a graph ignoring edge di-375 rections. Starting at an arbitrarily-chosen node. The *i*th node found by DFS must have been found from one of the 376 377 1 through i - 1 nodes.

378 2.4. Generalizing to more GNNs

Many highly performant GNN architectures can be perfectly represented by linear and quadratic constraints, and many more can be closely approximated. For example, if we choose our aggregation function to be a feature-wise average instead of a feature-wise sum, we can simply modify 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 decision variables $D^{(i)}$ for each layer, where row j of $D^{(i)}$ will represent the feature-wise average of the neighbors of node j. To properly constrain $D^{(i)}$, we add the constraint $1(\mathbf{1}^T A)D^{(i)} = AX$ to the model. The elements of D can be distributed in the multiplications for the left hand side, and then all the terms can be linearized as previously described. Now, constraint (3) can be changed to:

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

Consider a message passing layer from a Graph Isomor-379 380 phism Network [29], where updated node representations are calculated as $X' = h((A + (1 + \epsilon)I)X), h$ is a neural 381 network, and ϵ is a constant. We can split this computation 382 by constraining intermediate decision variables according 383 384 to the inner piece, $AX + ((1 + \epsilon)I)X$, and the application of the neural network to those intermediate variables, both 385 of which we previously discussed how to express with lin-386 ear constraints. Additionally, piecewise-linear approxima-387 tions can be created for non-linear functions, allowing us to 388 389 model different activation functions and the convolutional 390 layers in GCNs or GATs.

2.5. Optimality

We now show that a globally optimal solution to the MIP problem described above always exists, and that it can be provably found by MIPExplainer.

Theorem 1. Consider the MILP problem that maximizes 395 the objective function $X_c^{(L)} - y$, with decision variables A, 396 $X^{(i)}, \Phi^{(i)}, Z^{(i)}, B^{(i)}, E^{(i)}, d_j, y$ subject to constraints (4-15). A global optimum exists for this MILP. 398

Proof. Since all the constraints are linear equalities or in-399 equalities (after linearizing the multiplication of binary and 400 continuous variables in Eq.3), they define a polyhedral fea-401 sible region in the space parameterized by the decision vari-402 ables if binary variables are relaxed to be $\in [0, 1]$. All deci-403 sion variables in the MILP are bounded, either directly or in 404 terms of bounds on the input variables A and X, so the fea-405 sible set is a closed polytope (a compact set). The objective 406 function is linear in terms of the decision variables as cal-407 culated by $X_c^{(L)} - y$. Therefore, the linear relaxation of the 408 MILP must have an optimum, over the compact set, that is 409 located specifically on its boundary. In addition, there are a 410 finite number of integer solutions within the compact feasi-411 ble region, so at least one of them must have the maximum 412 objective value. 413 \square

We can use typical methods such as branch and bound, 414 cutting planes, and heuristics to find the optimal solution 415 efficiently. Using branch and bound, we start by finding 416 the optimal solution of the LP relaxation (i.e. the MIP with 417 the integrality constraints removed) of the MIP problem, for 418 example using the simplex method. This serves as an up-419 per bound to the original problem with integer domains for 420 some decision variables. If some integral decision variables 421 take fractional values in the relaxation, we branch on one of 422 them by partitioning the set of candidate solutions for the 423 problem with integer constraints into 2 subproblems, one 424 with the extra constraint that a chosen fractional variable is 425 at most the floor of its value in the LP relaxation's optimum 426 and another ensuring that the variable at least the ceiling 427 of that value. The optimal solution to the integral problem 428 will be the maximum optimal solution of these two subprob-429 lems, which can be solved recursively in the same way. On 430 the other hand, any integer solution serves as a lower bound 431 to the integral problem's optimum, and we can use heuris-432 tics to find increasingly better solutions as we descend in 433 the search tree and fix more variables. If the linear relax-434 ation solved at an internal node is infeasible or has a max-435 imum objective value that is lower than our current lower 436 bound, we have proven that the optimal solution to the in-437 tegral problem cannot lie anywhere in the subtree rooted at 438 that node, allowing us to skip all the nodes it contains in our 439 search. The process stops when there are no more subprob-440 lems to explore, at which point we will have found the opti-441 mal solution to the integral problem. In our experiments, we 442

517

use Gurobi Optimizer [7] to find the optimal solution efficiently. This solver benefits from being fully deterministic,
allowing us to generate explanations with a high degree of
consistency.

447 **2.6.** Practical Considerations

In practice, it can be difficult to solve MILPs correspond-448 449 ing to large GNNs, and several techniques are needed to 450 make the process tractable. Often, just finding an initial setting for all of the decision variables that satisfies all con-451 452 straints is difficult. In our experiments, we found that this step can actually take longer than the subsequent optimiza-453 tion. This problem can be completely eliminated with a 454 455 warm start. Starting from an arbitrary input graph (either from the dataset or not), we can compute a forward pass 456 through the network to obtain a valid setting of initial val-457 ues for almost all of the decision variables. In cases where 458 additional constraints have been imposed on the graph, such 459 as to ensure connectivity as described above, an input graph 460 must be converted into the canonical form that also satisfies 461 462 these constraints.

463 While a single, large number can be used to bound all of the continuous decision variables, tighter bounds greatly re-464 465 duce the time needed to compute optimal solutions. While automated bound-tightening procedures exist, it is faster to 466 use knowledge of the problem to bound manually. Each 467 hidden representation computed by the model is encoded 468 by a separate set of decision variables. Assuming we have 469 470 bounded the decision variables for one, we can compute 471 bounds for the outputs of a following transformation. For example, given a hidden representation vector x with lower 472 473 bound x_L and upper bound x_U , we can bound the output of a linear layer x' = Wx + b below and above by 474

475
$$x'_L = \operatorname{ReLU}(W)x_L + \operatorname{ReLU}(-W)x_U + b$$

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

477 Given bounds on the the decision variables representing the explanation graph, the input to the GNN, we can follow the 478 propagation of values through the GNN to iteratively bound 479 the set of decision variables for each hidden representation. 480 481 Bounds for the outputs of ReLU activation layers will be 482 the same as their inputs, but clipped below at 0. In the case 483 of layers like GraphSAGE convolutions where the output is the sum of several matrix multiplications, bounds can be 484 485 derived for each term in the sum and then added together.

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. 493 494 495 496 497 498

Despite these measures, runtime remains the most signif-499 icant drawback of the proposed approach. In our largest ex-500 periments, we were not able to guarantee a global optimum 501 within a single day. However, the largest reason for this run-502 time is the amount of symmetry in our formulation. A single 503 graph corresponds to a number of adjacency matrices that, 504 in the worst case, grows exponentially with the number of 505 vertices it contains. This hinders our ability to tighten the 506 upper bound while exploring the search tree, since an ex-507 isting global optimum may be transformed into another as 508 we traverse a branch. In the future, we plan to address this 509 problem by introducing additional constraints to reduce the 510 number of feasible adjacency matrices in each equivalence 511 class defined by graph isomorphism. Despite the increased 512 time needed to prove optimality, the proposed method often 513 finds optimal solutions early in the search. Therefore, we 514 impose a time limit during experiments to prove its practi-515 cality. 516

3. Empirical Evaluation

We use two synthetic datasets and one real-world dataset to 518 evaluate our method: Is_Acyclic, Shapes, and MUTAG. The 519 Is_Acyclic dataset comes from XGNN's experiments, and 520 has two classes consisting of cyclic and acyclic graphs of 521 various types. The cyclic graphs include graphs like grids, 522 single cycles, and wheels, while the acyclic class includes 523 graphs like paths and various types of trees. Every node is 524 given the same feature, a single constant, in order to isolate 525 the explanation methods' ability to capture structural infor-526 mation. For the Shapes dataset, which comes from GN-527 NInterpreter's experiments, graphs are first generated from 528 one of 5 base classes: lollipop graphs contain a fully con-529 nected component with one connection to a path graph's end 530 node, grid graphs are lattices where each internal node has 4 531 neighbors, star graphs have multiple outer nodes connected 532 to a single central node, and wheel graphs are star graphs 533 with a single cycle connecting the outer nodes. For each 534 of these graphs, a uniform proportion between 0 and 0.2 is 535 chosen, and the number of edges in the graph is increased by 536 that amount by adding in edges uniformly at random. The 537 features of each node are the same as in Is_Acyclic. The 538 MUTAG dataset [4] consists of graphs of chemical com-539 pounds, where nodes represent atoms and edges represent 540 bonds between them. Each compound is classified as be-541 ing either mutagenic or non-mutagenic. As described by 542 the creators of this dataset and in [9], mutagenic molecules 543 tend to have higher numbers of fused rings of carbon atoms. 544

(16)

CVPR 2024 Submission #16514. C	CONFIDENTIAL REVIEW	COPY. DO NOT DISTRIBUTE.
--------------------------------	---------------------	--------------------------

	# of Graphs	# of Classes	Average # of Nodes	Average # of Edges	# of Node Features
Is_Acyclic	533	2	28.5	68.1	1
Shapes	8000	5	27.2	144.9	1
MUTAG	188	2	17.9	39.6	7

Table 1. Dataset Summary

545 For this dataset, each node's features are a one-hot vector 546 indicating atom type.

547 **3.1.** Experimental Setup

548 Every dataset was randomly split into a training set (80%) and a test set (20%). GNNs were trained on each, and 549 550 performance metrics are reported in Table 2. In all experiments, the GNNs use GraphSAGE-style convolutions 551 552 with sum being used as the aggregation operator, followed 553 by a global mean pooling layer, and finally several fully-554 connected (FC) layers. ReLU activations are placed be-555 tween each hidden layer. For the Is_Acyclic and Shapes datasets, the GNN uses 2 convolutional layers computing 556 557 16 features per node, a FC layer computing 8 features, and 558 a final FC layer to compute the class logits. For the MUTAG dataset, the GNN uses 2 convolutional layers computing 64 559 560 and 32 features per node, two FC layers computing 16 and 8 561 features per graph, and a final FC layer to compute the log-562 its. We implemented these GNNs using PyTorch-Geometric [6]. Models were trained for 200 epochs, optimizing with 563 Adam [11] with a learning rate of 10^{-3} and L2 regulariza-564 tion with weight 10^{-4} . 565

For the experiments with XGNN, we used the imple-566 mentation provided by the authors in DIG^2 [16]. For the 567 568 experiments with GNNInterpreter, we also use the implementation provided by the authors³. For the MUTAG 569 570 dataset, XGNN's graph generator policy network was pe-571 nalized when it violated valence constraints while gener-572 ating molecules, and no penalties were used on the other 573 datasets. Both methods optimize objective functions with 574 weighted regularizers, but there is no single quantitative 575 metric for model-level explanation quality, so it is not possible to objectively tune the hyperparameters associated with 576 577 these methods. Instead, we used default sets of hyperparameters provided in the implementations. An exception was 578 579 made for XGNN, the default regularization weights caused 580 the graph generator to quickly learn a policy that stopped after the first node in several instances. To fix this, we in-581 582 creased the reward for creating additional valid edges until it became favorable for the model to generate reasonably-583 584 sized explanations.

In the experiments with MIPExplainer, adjacency ma trices were constrained to be symmetric to represent undi rected connected graphs without self-loops. For the MU-

	Train Accuracy	Test Accuracy	# of Model Parameters
Is_Acyclic	0.998	1.000	730
Shapes	0.991	0.993	757
MUTAG	0.893	0.895	5770

Table 2. Performance Metrics of Trained GNNs

TAG dataset, node features were constrained to one-hot vectors by ensuring the sum of the elements in each row added588up to 1. Any experiments lasting longer than 6 hours were590automatically terminated, and we report the best solution591found. To ensure that any resemblance to target classes592would not come from an initial solution, all runs for our593method were initialized with a path graph of n nodes.594

3.2. Results

The main results from our experiments are shown in Table 596 3. Note that when depicting molecular graphs, the node 597 colors are assigned as follows: gray=C, blue=N, red=O, 598 cyan=F, purple=I, green=Cl, and brown=Br. In the exper-599 iments with Is_Acyclic, MIPExplainer explains the cyclic 600 class with a fully-connected graph, which has the maximum 601 possible number of cycles, and the acyclic class with a star 602 graph, which is one of the most straightforward examples 603 from the class. XGNN does explain the acyclic class with 604 a graph that has fewer cycles than the explanation for the 605 cyclic class, but both graphs are connected somewhat ran-606 domly and both contain cycles. On the other hand, GN-607 NInterpreter produces a more densely-connected graph to 608 explain the acyclic class than the cyclic class. For MIPEx-609 plainer, the solver was able to prove the optimality of both 610 classes, taking 1.52 seconds for the cyclic class explanation 611 and 227.68 seconds for the acyclic class explanation. The 612 left plot in Figure 4 shows how the bounds converged over 613 the course of the latter experiment. This demonstrates how 614 graph symmetries factor into MIPExplainer's runtime. A 615 fully connected graph with equal node features only has a 616 single adjacency matrix and feature matrix representation, 617 while a star graph with n nodes has n representations, as 618 there are n options for the position of the central node in 619 the node ordering. As a result, more of the search tree must 620 be explored to prove optimality. The right plot in Figure 4 621 shows the number of explored nodes and unexplored leaf 622 nodes over the course of the search for the optimal solution. 623 A total of 183543 nodes were explored to verify the optimal 624 solution. 625

For the Shapes dataset, we can easily recognize the 626 classes of each of the explanations generated by MIPEx-627 plainer. Despite the fact that a significant amount of noise 628 was added to the training data, the explanations are rela-629 tively clean. For reference, three examples from the wheel 630 class are shown in Figure 3. In the cases of lollipops and 631 stars, we see important features of the graph duplicated, a 632 lollipop with two ends and star with two centers. Because 633 we chose a higher number of nodes to make the patterns 634

²https://github.com/divelab/DIG

 $^{^{3}\}mbox{https}$: / / github . com / yolandalalala / GNNInterpreter



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



Figure 2. Smaller explanations for lollipops, wheels, grids, and stars respectively



Figure 3. Three randomly selected wheel graphs from the Shapes dataset

clearer, optimality was not proven for these explanations,
but they still appear to be reasonable. To further demonstrate that MIPExplainer performs well at varying choices
of the number of nodes in the explanation, we provide 5node explanations for the same classes in Figure 2, which
were proven to be optimal.

For the MUTAG dataset, our results show that 641 the model's true behavior may not accurately describe 642 643 chemistry-related patterns in the data. For the mutagenic class, the produced explanation is simply a complete graph 644 645 of carbon atoms. Comparing this with the result from the cyclic class of the cyclicity dataset, it seems that the model 646 may simply be looking for arbitrary cycles of carbon atoms, 647 not just ones that occur in molecules. Neither of the expla-648 nations generated by the two baseline methods contained 649 650 a cycle of carbon atoms. The explanation of the non-651 mutagenic class are not as reasonable across all methods, which is expected since non-mutagens are more accurately 652 653 described by the absence of mutagenic features than by nonmutagenic features. The generated explanation mostly con-654 655 sists of bromine atoms, which only actually appear in 2 of 656 the graphs in the dataset. Despite the larger network ar-657 chitecture, both these solutions were able to be verified as optimal within the time limit. 658

4. Conclusion and Discussion

660 Despite the ability of GNNs to model complex patterns in
661 graph-structured data, their lack of transparency remains
662 one of the key factors hindering their application in a wide



Figure 4. Solver metrics while explaining the acylcic class of Is_Acyclic: Bounds converging to the globally optimum (left), and the number of explored/unexplored nodes in the search tree (right)

range of domains. Model-level explanations of these net-663 works are key to understanding the information they learn 664 and improving their trust and reliability. In order to ad-665 dress key shortcomings that limit the use of existing meth-666 ods in most real-world situations, this work proposes MIP-667 Explainer for generating such explanations. Without a way 668 to objectively evaluate their quality, it is essential that gen-669 erated explanations are truly high-quality solutions of opti-670 mization problems that are not sensitive to user-defined hy-671 perparameters. MIPExplainer achieves this by avoiding the 672 use of both weighted regularizers and stochastic optimiza-673 tion, instead focusing on maximizing a simpler objective 674 with deterministic methods that are able to prove the global 675 optimality of the generated solutions. Minimal assumptions 676 are made about the distributions of graphs and their features, 677 and no secondary models are trained in the process. 678

The proposed method also has several shortcomings, 679 which we hope to address in future work. While it is more 680 general than previous methods in some ways, it also re-681 quires different GNN layers to be individually encoded with 682 constraints, and may require piecewise-linear approxima-683 tions for highly nonlinear components. From a practical 684 perspective, the runtime of MIPExplainer as described here 685 is the most significant. Reducing symmetries in the encod-686 ing can greatly improve runtime, but this is a hard problem 687 in general, and more work is required to understand which 688 symmetries are the most costly when optimizing over sets 689 of graphs in this way. Despite these limitations, even before 690 proving optimality, we observe that the proposed method is 691 generally able to find reasonable explanations. 692

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

777

778

779

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799

801

802

803

804

805

806

693 References

- [1] Navid Ansari, Hans-Peter Seidel, and Vahid Babaei. Mixed
 integer neural inverse design. ACM Transactions on Graph ics, 41(4):151:1–151:14, 2022. 2
- 697 [2] Steve Azzolin, Antonio Longa, Pietro Barbiero, Pietro Liò,
 698 and Andrea Passerini. Global Explainability of GNNs via
 699 Logic Combination of Learned Concepts. 2022. 1, 2
- [3] Myun-Seok Cheon. An outer-approximation guided optimization approach for constrained neural network inverse problems. *Mathematical Programming: Series A and B*, 196 (1-2):173–202, 2022. 2
- [4] Asim Kumar Debnath, Rosa L. Lopez De Compadre, Gargi Debnath, Alan J. Shusterman, and Corwin Hansch.
 Structure-activity relationship of mutagenic aromatic and heteroaromatic nitro compounds. Correlation with molecular orbital energies and hydrophobicity. *Journal of Medicinal Chemistry*, 34(2):786–797, 1991. 6
 - [5] Alexandre Duval and Fragkiskos D. Malliaros. GraphSVX: Shapley Value Explanations for Graph Neural Networks. In Machine Learning and Knowledge Discovery in Databases. Research Track, pages 302–318, Cham, 2021. Springer International Publishing. 1
 - [6] Matthias Fey and Jan E. Lenssen. Fast Graph Representation Learning with PyTorch Geometric. In *ICLR Workshop on Representation Learning on Graphs and Manifolds*, 2019. 7
 - [7] Gurobi Optimization, LLC. Gurobi Optimizer Reference Manual, 2023. 6
 - [8] William L. Hamilton, Rex Ying, and Jure Leskovec. Inductive Representation Learning on Large Graphs. 2017. 1
 - [9] Kuo-Hsiang Hsu, Bo-Han Su, Yi-Shu Tu, Olivia A. Lin, and Yufeng J. Tseng. Mutagenicity in a Molecule: Identification of Core Structural Features of Mutagenicity Using a Scaffold Analysis. *PLOS ONE*, 11(2):e0148900, 2016. 6
 - [10] Erwin Kalvelagen. Multiplication of a continuous and a binary variable, 2008. 4
 - [11] Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. In 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings, 2015. 7
 - [12] Thomas N. Kipf and Max Welling. Semi-Supervised Classification with Graph Convolutional Networks. 2016. 1
- [13] Wanyu Lin, Hao Lan, and Baochun Li. Generative causal explanations for graph neural networks. In *International Conference on Machine Learning*, pages 6666–6679. PMLR, 2021. 1
- [14] Wanyu Lin, Hao Lan, and Baochun Li. Generative Causal
 Explanations for Graph Neural Networks. In *International Conference on Machine Learning*, 2021. 2
- [15] Wanyu Lin, Hao Lan, Hao Wang, and Baochun Li. OrphicX:
 A Causality-Inspired Latent Variable Model for Interpreting
 Graph Neural Networks. 2022. 2
- [16] Meng Liu, Youzhi Luo, Limei Wang, Yaochen Xie, Hao
 Yuan, Shurui Gui, Zhao Xu, Haiyang Yu, Jingtun Zhang, Yi
 Liu, Keqiang Yan, Bora Oztekin, Haoran Liu, Xuan Zhang,
 Cong Fu, and Shuiwang Ji. DIG: A Turnkey Library for
 Diving into Graph Deep Learning Research. *arXiv preprint arXiv:2103.12608*, 2021. 7

- [17] Ana Lucic, Maartje A Ter Hoeve, Gabriele Tolomei, Maarten De Rijke, and Fabrizio Silvestri. Cf-gnnexplainer: Counterfactual explanations for graph neural networks. In *International Conference on Artificial Intelligence and Statistics*, pages 4499–4511. PMLR, 2022. 1
 754
- [18] Dongsheng Luo, Wei Cheng, Dongkuan Xu, Wenchao Yu, Bo Zong, Haifeng Chen, and Xiang Zhang. Parameterized explainer for graph neural network. *Advances in neural information processing systems*, 33:19620–19631, 2020. 1
- [19] Alan Perotti, Paolo Bajardi, Francesco Bonchi, and André Panisson. Explaining Identity-aware Graph Classifiers through the Language of Motifs. In 2023 International Joint Conference on Neural Networks (IJCNN), pages 1–8, 2023. ISSN: 2161-4407. 1
- [20] Phillip E Pope, Soheil Kolouri, Mohammad Rostami, Charles E Martin, and Heiko Hoffmann. Explainability methods for graph convolutional neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10772–10781, 2019. 1
- [21] Sayan Saha, Monidipa Das, and Sanghamitra Bandyopadhyay. GraphEx: A User-Centric Model-Level Explainer for Graph Neural Networks. In *The First Tiny Papers Track at ICLR 2023, Tiny Papers @ ICLR 2023, Kigali, Rwanda, May 5, 2023.* OpenReview.net, 2023. 2
- [22] Sayan Saha, Monidipa Das, and Sanghamitra Bandyopadhyay. GraphEx: A User-Centric Model-Level Explainer for Graph Neural Networks. 2023. 2
- [23] Michael Sejr Schlichtkrull, Nicola De Cao, and Ivan Titov. Interpreting graph neural networks for nlp with differentiable edge masking. arXiv preprint arXiv:2010.00577, Proceedings of International Conference on Learning Representations, 2020. 1
- [24] Thomas Schnake, Oliver Eberle, Jonas Lederer, Shinichi Nakajima, Kristof T Schütt, Klaus-Robert Müller, and Grégoire Montavon. Higher-order explanations of graph neural networks via relevant walks. *IEEE transactions* on pattern analysis and machine intelligence, 44(11):7581– 7596, 2021. 1
- [25] Yong-Min Shin, Sun-Woo Kim, and Won-Yong Shin. PAGE: Prototype-Based Model-Level Explanations for Graph Neural Networks. 2022. 2
- [26] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph Attention Networks. 2017. 1
- [27] Minh Vu and My T Thai. Pgm-explainer: Probabilistic graphical model explanations for graph neural networks. Advances in neural information processing systems, 33:12225– 12235, 2020. 1
- [28] Xiaoqi Wang and Han-Wei Shen. GNNInterpreter: A Probabilistic Generative Model-Level Explanation for Graph Neural Networks. 2022. 2, 3
- [29] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka.How Powerful are Graph Neural Networks? 2018. 5
- [30] Zhitao Ying, Dylan Bourgeois, Jiaxuan You, Marinka Zitnik, and Jure Leskovec. Gnnexplainer: Generating explanations for graph neural networks. *Advances in neural information processing systems*, 32, 2019. 1

- 807 [31] Junchi Yu, Tingyang Xu, Yu Rong, Yatao Bian, Junzhou
 808 Huang, and Ran He. Graph Information Bottleneck for Sub809 graph Recognition. 2020. 1
- [32] Hao Yuan, Jiliang Tang, Xia Hu, and Shuiwang Ji. XGNN:
 Towards Model-Level Explanations of Graph Neural Networks. In *Proceedings of the 26th ACM SIGKDD Interna- tional Conference on Knowledge Discovery & Data Mining*,
 pages 430–438, Virtual Event CA USA, 2020. ACM. 2
- [33] Hao Yuan, Haiyang Yu, Jie Wang, Kang Li, and Shuiwang Ji.
 On explainability of graph neural networks via subgraph explorations. In *International conference on machine learning*, pages 12241–12252. PMLR, 2021. 1
- [34] Hao Yuan, Haiyang Yu, Shurui Gui, and Shuiwang Ji. Explainability in Graph Neural Networks: A Taxonomic Survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(5):5782–5799, 2023. 1
- [35] Zaixi Zhang, Qi Liu, Hao Wang, Chengqiang Lu, and Cheekong Lee. ProtGNN: Towards Self-Explaining Graph Neural Networks. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(8):9127–9135, 2022. 1