

Providing Bipartite GNN Explanations with PGExplainer

Tristan Marten Lewin Schulz

Bachelor of Science

May 19, 2025

Supervisors:

Prof. Dr. Stefan Harmeling

Lukas Schneider

Artificial Intelligence (VIII)

Department of Computer Science

TU Dortmund University

Abstract

Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua.

Contents

1	Introduction	1
2	Related Work	3
3	Theoretical Background	5
3.1	Deep learning	5
3.1.1	Multilayer Perceptron	6
3.1.2	Backpropagation	7
3.1.3	TODO: Regularization	7
3.1.4	Batch Normalization	7
3.1.5	Monte Carlo Sampling	8
3.2	Graph Theory	8
3.3	Information Theory	10
3.3.1	Entropy	10
3.3.2	Relative Entropy and Cross-Entropy	11
3.3.3	Mutual Information	11
3.4	Graph Neural Networks	11
3.5	Perturbation-based Explainability in GNNs	14
3.6	Boolean Satisfiability Problem	15
3.6.1	Representation as Bipartite Graph	16
3.6.2	Unsatisfiable Cores	17
4	PGExplainer - Main part	19
4.1	Theory	20
4.1.1	Learning Objective	20
4.1.2	Reparameterization Trick	21
4.1.3	Global Explanations	22
4.1.4	Regularization Terms	24
4.2	Reimplementation	27
4.3	Application on NeuroSAT	28
5	Results	31

6 Discussion	33
7 Conclusion	35
Bibliography	36
A Supplementary Material	41
Affidavit	42

Chapter 1

Introduction

SAT motivation: Advancements in deep learning SAT solving, e.g. NeuroSAT. Need for evaluation of these models! use of GNN explainers, since SAT can be reduced to graph domain. GOAL: Application of PGExplainer on NeuroSAT to generate explanations. Explanations need gt to be evaluated on accuracy. Use concepts like unsat cores and backbones as gt and compare to explanations provided by PGExplainer to see whether NeuroSAT explanations align with "human-observable" principles. Concrete: Graph task: Prediction of UNSAT and unsat cores as gt. Node task: Difficult with NeuroSAT?

Therefore explanation and replication of PGExplainer. After successful replication, application on NeuroSAT.

PGExplainer claims to be model-agnostic - Works for any GNN

Since our goal is application on a bipartite problem, reimplementing of PGExplainer that changes architecture of target GNN to a PyG layer that works on bipartite graphs and also allows passing of edge weights (required for PGE)

Pytorch reimplementations already exist, eg. Replication study RE-PGE that focuses on original code to replicate results. We follow paper as close as possible, and also conduct the code to try different settings/hyperparameters for our changed GNN - OUR FOCUS: Evaluate whether PGExplainer still applicable for our slightly changed GNN model!

Therefore, use PGExplainer and Replication paper as baselines to compare our results

Lastly, apply PGExplainer NeuroSAT, a solver for bipartite graph problem SAT

Chapter 2

Related Work

TODO: PGExplainer

Original work by Luo et al.[17]; Original PGExplainer? SHORT SUMMARY. We seek to provide a reimplementation using PyG framework and reproduce the results, while also applying the framework to another domain, namely SAT problem in the form of bipartite graphs.

TODO: GNNExplainer

Baseline for PGExplainer by Yuan et al.[28]; DIFFERENCE TO PGEXPLAINER

TODO: RE-PGExplainer

Replication study of PGExplainer using pytorch by Holdijk et al.[12]; Unable to achieve results of original implementation. Improvement on some datasets, considerable deterioration on BA-2Motif. Criticize approach of original and lacklustre documentation as well as differences between codebase and paper. We try own reimplementation that follows the original paper, as well as code and reimplementation for uncertainties + use a slightly different architecture in underlying GNN + Hyperparameter search of combination of parameters used in original and replication. Treated as additional baseline

TODO: Taxonomy paper? Taxonomy paper evaluates PGExplainer on Fidelity and achieves low scores, sides with reproducibility paper above: "is not performing as promising as its original reported results". Propose only using PGExplainer for Node classification task.

NeuroSAT

NeuroSAT by Selsam et al.[25] is a machine learning approach for solving the boolean satisfiability problem using a message passing neural network. It is able to detect satisfying assignments but lacks proofs of unsatisfiability. The authors performed a small study on the detection of unsat cores, revealing that NeuroUNSAT is able to detect UNSAT cores if the UNSAT problems contain a specific UNSAT core. However, this is expected to be due to the model memorizing the unsat cores, rather than generalizing to any unsat core.

We want to test this by applying the PGExplainer to the NeuroSAT model, and evaluate whether the generated explanations do align with unsat cores. TODO: Code provided by rodhi?

Chapter 3

Theoretical Background

In this chapter we define the necessary background for understanding PGExplainer as well as the follow-up work regarding its application on the Boolean Satisfiability Problem (SAT).

3.1 Deep learning

In this chapter we introduce Deep Learning (DL) in the context of Machine Learning (ML) and their concepts required for this work. The definitions in this chapter loosely follow Goodfellow et al.[10].

An ML algorithm generally learns to perform a certain task from data. Common ML tasks include classification, where the goal is to assign an input to one of k categories, and regression, where the program shall predict a numeric value given some input.

ML algorithms can be broadly divided into supervised and unsupervised learning. In this work, we focus on supervised learning, meaning that the algorithm learns from a dataset containing both features and labels or targets that the algorithm is supposed to predict. In other words, the algorithm learns from a training set of inputs x and outputs y to associate another unseen input with some output.

DL entails expanding the size of the model used in our algorithm to allow for representations of functions with increasing complexity. This enables many human-solvable tasks that consist of mapping an input vector to an output vector to be performed with DL, given sufficiently large models and datasets.

In many cases DL involves optimizing a function, usually by minimizing $f(x)$. This objective function is also referred to as loss function in the context of minimization. To minimize a function $f(x)$ we make use of the derivative $f'(x)$ that tells us how to change x in order to get an improvement in y . We can reduce $f(x)$ by moving x in the direction opposite of the sign of its derivative, since $f(x - \epsilon \text{sign}(f'(x))) < f(x)$ for small enough ϵ . This technique is called gradient descent (see Cauchy[3]).

Furthermore, a DL algorithm typically consists of the following components: a dataset specification, an objective function, an optimization procedure, like gradient descent, and a model.

3.1.1 Multilayer Perceptron

A classical DL model is the multilayer perceptron (MLP) with the general goal of approximating a function f^* . In the case of classification we could define a function $y = f^*(x)$ that maps an input x to a label y . The MLP then defines the mapping $y = f(x; \Theta)$ and learns the value of the parameters Θ that best approximate the function. These models are also referred to as feedforward neural networks, as they process information from x , through the intermediate computations that define f , to the output y without feedback connections that would feed outputs back to itself. The name network is derived from their representation as a composition of multiple different functions, that are described by a directed acyclic graph. An example network is $f(x) = f^{(3)}(f^{(2)}(f^{(1)}(x)))$ that consists of input layer $f^{(1)}$, hidden layer $f^{(2)}$ and output layer $f^{(3)}$. The length of this chain of functions is called depth and origin of the term "deep learning". The approximation is achieved by training our network with training data, that consists of approximated examples of $f^*(x)$ at different points in the training and labels $y \approx f^*(x)$. These training examples dictate the output layer to generate a value close to y for each x . The learning algorithm then learns to utilize the other hidden layers, without specified behaviors, to achieve the best approximation. It is to note that the hidden layers are vector-valued with each vector element, referred to as unit, loosely taking the role of a neuron in neuroscience. Models are therefore also referred to as neural networks.

A linear layer for our model with parameters Θ consisting of weight w and bias b can be described as $f(x; w, b) = x^T w + b$ for an input vector x . To keep the model from strictly learning a linear function of its inputs we can calculate the values of a layer h by applying an activation function g to its output to describe the features. This results in our hidden layer $h = g(W^T x + c)$, with W containing the weights of a linear transformation and c the biases.

The activation function usually serves the purpose of mapping to a real number between 0 and 1, imitating the activation of a neuron. We try to use activation functions that are continuous differentiable and easily calculated, to minimize computational complexity (see Liu et al.[16]). Examples are the Sigmoid function

$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad (3.1)$$

and the Rectified Linear Unit (ReLU)

$$ReLU(x) = \begin{cases} 0 & x \leq 0, \\ 1 & x > 0. \end{cases} \quad (3.2)$$

3.1.2 Backpropagation

The back propagation algorithm is commonly used during neural network training. It optimizes the network parameters by leveraging gradient descent. It first calculates the values for each unit in the network given the input and a set of parameters in a forward order. Then the error for each variable to be optimized is calculated, and the parameters are updated according to their corresponding partial derivative backwards. These two steps will repeat until reaching the optimization target.

TODO: Concrete formulas with chain rule of calculus?

3.1.3 TODO: Regularization

"A central problem in machine learning is how to make an algorithm that will perform well not just on the training data, but also on new inputs. Many strategies used in machine learning are explicitly designed to reduce the test error, possibly at the expense of increased training error. These strategies are known collectively as regularization."

Parameter Norm Penalties (e.g. Weight decay)

Early stopping "When training large models with sufficient representational capacity to overfit the task, we often observe that training error decreases steadily over time, but validation set error begins to rise again." "This means we can obtain a model with better validation set error (and thus, hopefully better test set error) by returning to the parameter setting at the point in time with the lowest validation set error." "Every time the error on the validation set improves, we store a copy of the model parameters. When the training algorithm terminates, we return these parameters, rather than the latest parameters." "algorithm terminates when no parameters have improved over the best recorded validation error for some pre-specified number of iterations. "

Dropout

3.1.4 Batch Normalization

In deep models composed of several layers the gradient tells us how to update each parameter assuming that the other layers do not change. Since in practice all layers are updated simultaneously, this may lead to unexpected results. Batch normalization is an optimization technique that can be applied to any input or hidden layer to reduce the aforementioned problem[10].

Let \mathbf{B} be a minibatch of activations of the layer to normalize, in the form of a design matrix with rows of activations. To normalize \mathbf{B} , it is replaced by

$$\mathbf{B}' = \frac{\mathbf{B} - \boldsymbol{\mu}}{\boldsymbol{\sigma}}, \quad (3.3)$$

where $\boldsymbol{\mu}$ is a vector containing the mean of each unit and $\boldsymbol{\sigma}$ is a vector containing the standard deviation of each unit. Each activation is normalized individually with the $\boldsymbol{\mu}$

and σ corresponding to its row. The network processes \mathbf{B}' as functionally equivalent to \mathbf{B} . During training,

$$\mu = \frac{1}{m} \sum_i \mathbf{B}_i \quad (3.4)$$

and

$$\sigma = \sqrt{\delta + \frac{1}{m} \sum_i (\mathbf{B}_i - \mu)^2}, \quad (3.5)$$

where δ is a small positive value introduced to avoid \sqrt{z} at $z = 0$ and m is the number of elements in the batch. It is important to note that back propagation is performed through all three operations, to prevent the gradient from proposing an operation that solely increases the mean and standard deviation of the output of a layer. Batch normalization therefore reparameterizes the model to include some units that are standardized by definition.

When testing the model, we replace μ and σ with running averages that were collected during training, to enable the evaluation of individual examples outside minibatches.

3.1.5 Monte Carlo Sampling

In order to best approximate a randomized algorithm we can make use of Monte Carlo methods as described in Goodfellow et al.[10][p.590]. A common practice in ML is to draw samples from a probability distribution and using these to form a Monte Carlo estimate of some quantity. This can be used to train a model that can then sample from a probability distribution itself.

More specifically, the idea of Monte Carlo sampling is to view a sum as if it was an expectation under some distribution and to approximate this estimate with a corresponding average. Let

$$s = \sum_x p(x) f(x) = E_p[f(x)] \quad (3.6)$$

be the sum to estimate with p being a probability distribution over a random variable x . Then s can be approximated by drawing n samples from p and constructing the empirical average

$$\hat{s}_n = \frac{1}{n} \sum_{i=1}^n f(x^{(i)}). \quad (3.7)$$

3.2 Graph Theory

These definitions will loosely follow Liu et al.[16]. A graph is a data structure consisting of a set of nodes that are connected via edges, modeling objects and their relationships. It can be represented as $G = (V, E)$ with $V = \{v_1, v_2 \dots v_n\}$ being the set of n nodes, and $E \in V \times V$ the set of edges. We adopt the conventions from Diestel[7][p.2] to refer to the node and edges set of any graph G with $V(G)$ and $E(G)$ respectively, regardless of the actual names of the sets, as well as a referring to G with node set V as G on V . An

edge $e = (u, v)$, sometimes accompanied by an edge weight $e_{u,v} \in (0, 1)$, connects nodes u and v , making them neighbors. The neighborhood $\mathcal{N}(u)$ of node u in $V(G)$ is defined as $\mathcal{N}(u) = \{v \in V(G) \mid (v, u) \in E(G)\}$. We define the k -hop neighborhood of node u in $V(G)$ as $\mathcal{N}_k(u) = \{v \in V(G) \mid \text{dis}_G(v, u) \leq k\}$, where $\text{dis}_G(v, u)$ denotes the distance of nodes v, u in G , defined as the length of the shortest path between the two nodes. Edges are either directed or undirected and lead to directed or undirected graphs if exclusively present. The degree of a node v is the number of edges connected to v and denoted by $d(v)$. G can be described by an adjacency matrix $A \in \mathbb{R}^{n \times n}$, where Nodes in V are associated with the d -dimensional features, denoted by $X \in \mathbb{R}^{n \times d}$

$$A_{ij} = \begin{cases} 1 & \text{if } \{v_i, v_j\} \in E \text{ and } i \neq j, \\ 0 & \text{otherwise.} \end{cases}$$

If G is an undirected Graph the adjacency matrix will be symmetrical.

Alternatively an undirected graph $G = (V, E)$ with n nodes and m edges can be represented as an incidence matrix $M \in \mathbb{R}^{n \times m}$, where

$$M_{ij} = \begin{cases} 1 & \text{if } \exists k \text{ s.t. } e_j = \{v_i, v_k\}, \\ 0 & \text{otherwise.} \end{cases}$$

The diagonal degree matrix $D \in \mathbb{R}^{n \times n}$ of G is defined via:

$$D_{ii} = d(v_i).$$

G is called a subgraph of another graph $G' = (V', E')$ if $V(G) \subseteq V(G')$ and $E(G) \subseteq E(G')$. This is denoted as $H \subseteq G$. The number of nodes in a graph $|V|$ is its order and the number of edges $|E|$ is its size (see Diestel[7]). We additionally define bipartite graphs according to Asratian et al.[1]: A graph G is bipartite if the set of nodes V can be partitioned into two sets V_1 and V_2 so that no two nodes from the same set are adjacent. The sets V_1 and V_2 are called colour classes and (V_1, V_2) is a bipartition of G . This means that if a graph is bipartite all nodes in V can be coloured by at most two colours so that no two adjacent nodes share the same colour.

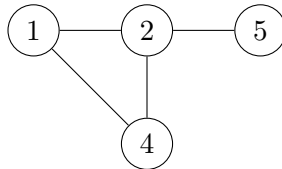


Figure 3.1: A simple undirected graph G with $V = \{1, 2, 4, 5\}$ and $E = \{\{1, 2\}, \{2, 4\}, \{1, 4\}, \{2, 5\}\}$.

Random Graphs

Gilbert[9] describes the process of generating a random graph of order N by assigning a common probability of existence to each potential edge between any two nodes. For each of these potential edges an experiment is performed independently to determine whether it shall be included in the resulting graph. Note that this process can be modeled using a Bernoulli distribution.

A random graph is further described by Diestel[7][p.323] as follows. Let $V = \{0, \dots, n-1\}$ be a fixed set of n elements. Say we want to define the set \mathcal{G} of all graphs on V as a probability space, which allows us to ask whether a Graph $G \in \mathcal{G}$ has a certain property. To generate our random graph we then decide from some random experiment whether e shall be an edge of G for each potential $e \in V \times V$. The probability of success - accepting e as edge in G - is defined as $p \in [0, 1]$ for each experiment. This leads to the probability of G being a particular graph G_0 on V with e.g. m edges being equal to $p^m q^{\binom{n}{2}-m}$ with $q := 1 - p$.

3.3 Information Theory

To fully understand the learning objective of PGExplainer it is necessary to define the concepts of entropy and mutual information. We follow the definitions by Cover et al.[5][p.13] if not stated otherwise.

3.3.1 Entropy

Entropy is used to describe the uncertainty of a random variable. It measures the amount of information required on average to describe a random variable. Let X be a discrete random variable with alphabet \mathcal{X} and probability mass function $p(x) = \Pr\{X = x\}$ for $x \in \mathcal{X}$. The entropy $H(X)$, also written as $H(p)$, is defined as

$$H(X) = - \sum_{x \in \mathcal{X}} p(x) \log p(x). \quad (3.8)$$

The log is to the base e and entropy is measured in nats in our case. TODO: DEFINE EDGE CASES $\log 0$

The conditional entropy of Y given X is defined as the expected value of the entropies of the conditional distributions, averaged over the conditioning random variable. If $(X, Y) \sim p(x, y)$ for a pair of discrete random variables (X, Y) with joint distribution $p(x, y)$, the conditional entropy is defined as

$$H(Y|X) = - \sum_{x \in \mathcal{X}} p(x) H(Y|X = x) \quad (3.9)$$

$$= - \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log p(y|x) \quad (3.10)$$

$$= -\mathbb{E} \log p(Y|X), \quad (3.11)$$

with expectation \mathbb{E} .

3.3.2 Relative Entropy and Cross-Entropy

The relative entropy between two distributions is a measure of "distance" between the two. It measures the inefficiency of assuming a distribution to be q when the true distribution is p . It is not a true measure of distance as it is among other things not symmetrical. The relative entropy takes a value of 0 only if $p = q$. We define the KL divergence or relative entropy between two probability mass functions $p(x), q(x)$ as

$$D_{KL}(p||q) = \sum_{x \in \mathcal{X}} p(x) \log \frac{p(x)}{q(x)} \quad (3.12)$$

Suppose we know the true distribution p of our random variable. We could then construct a code with an average description length of $H(p)$. If we used the code for the distribution q instead, we would need $H(p) + D_{KL}(p||q)$ nats to describe the random variable on average. This is also referred to as the cross-entropy (see Goodfellow et al.[10][p.74]):

$$H(p, q) = H(p) + D_{KL}(p||q) \quad (3.13)$$

We derive for the discrete case with mass probability functions p, q defined on the same support \mathcal{X} :

$$H(p, q) = H(p) + D_{KL}(p||q) = - \sum_{x \in \mathcal{X}} p(x) \log p(x) + \sum_{x \in \mathcal{X}} p(x) \log \frac{p(x)}{q(x)} \quad (3.14)$$

$$= - \sum_{x \in \mathcal{X}} p(x) \log p(x) + \sum_{x \in \mathcal{X}} p(x) \log p(x) - \sum_{x \in \mathcal{X}} p(x) \log q(x) \quad (3.15)$$

$$= - \sum_{x \in \mathcal{X}} p(x) \log q(x) \quad (3.16)$$

3.3.3 Mutual Information

Another closely related concept is mutual information. It measures the amount of information that one random variable contains about another or the reduction in uncertainty of said variable due to knowing the other. A high mutual information therefore implies that the information of one variable can be gathered from the other.

Let X and Y be two random variables with the joint probability mass function $p(x, y)$ and marginal probability mass functions $p(x)$ and $p(y)$. Mutual information $I(X; Y)$ is the relative entropy between the joint distribution and the product distribution $p(x)p(y)$:

$$I(X; Y) = \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log \frac{p(x, y)}{p(x)p(y)} \quad (3.17)$$

$$= H(X) - H(X|Y) \quad (3.18)$$

3.4 Graph Neural Networks

TODO: DOES THIS CONVEY MESSAGE PASSING???

Graph Neural Networks(GNNs)[23] are a DL-based approach that operates on graphs. Due to their unique non-Euclidean property, they find usage in classification, link prediction, and clustering tasks. Their high interpretability and strong performance have led to GNNs becoming a commonly employed method in graph analysis. They combine the key features of convolutional neural networks[15], such as local connection, shared weights, and multi-layer usage, with the concept of graph embeddings[2] (TODO: NON-EXISTENT?) to leverage the power of feature extraction and representation as low-dimensional vectors for graphs (see Liu et al.[16]).

Graphs are a common way of representing data in many different fields, including ML. ML applications on graphs can mostly be divided into graph-focused tasks and node-focused tasks. For graph-focused applications our model does not consider specific singular nodes, but rather implements a classifier on complete graphs. In node-focused applications however the model is dependent on specific nodes, leading to classification tasks that rely on the properties of each node. The supervised GNN model by Scarselli et a.[23] tries to preserve the important, structural information of graphs by encoding their topological relationships among nodes.

A node is naturally defined by its features as well as its related nodes in the graph. The goal of a GNN is to learn state embeddings $\mathbf{h}_v \in \mathbb{R}^S$ for each node v , that map the neighborhood of a node into a representation. These embeddings are used to obtain outputs \mathbf{o}_v , that e.g. may contain the distribution of a predicted node label. The GNN model proposed by Scarselli et al.[23] uses undirected homogeneous graphs with \mathbf{x}_v describing the d -dimensional features of each node and x_e the optional features of each edge. $co[v]$ and $ne[v]$ denote the set of edges and neighbors of node v respectively. The model updates the node states according to the input neighborhood with a local transition function f that is shared by all nodes. Additionally, the local output function g is used to produce the output of each node. \mathbf{h}_v and \mathbf{o}_v are therefore defined as

$$\mathbf{h}_v = f(\mathbf{x}_v, \mathbf{x}_{co[v]}, \mathbf{h}_{ne[v]}, \mathbf{x}_{ne[v]}), \quad (3.19)$$

$$\mathbf{o}_v = g(\mathbf{h}_v, \mathbf{x}_v), \quad (3.20)$$

with \mathbf{x} denoting input features and \mathbf{h} the hidden state. $\mathbf{x}_v, \mathbf{x}_{co[v]}, \mathbf{h}_{ne[v]}, \mathbf{x}_{ne[v]}$ denote the features of the node v and of its edges, as well as the states and features of its neighboring nodes, respectively. We define $\mathbf{H}, \mathbf{O}, \mathbf{X}$ and \mathbf{X}_N as the matrices that are constructed by stacking all states, outputs, features, and node features, respectively. This allows us to define with the global transition function F and the global output function G , which are stacked versions of their local equivalent for all nodes in a graph:

$$\mathbf{H} = F(\mathbf{H}, \mathbf{X}), \quad (3.21)$$

$$\mathbf{O} = G(\mathbf{H}, \mathbf{X}_N). \quad (3.22)$$

Note that F is assumed to be a contraction map and the value of \mathbf{H} is the fixed point of equation (3.21). To compute the state the iterative scheme

$$\mathbf{H}^{t+1} = F(\mathbf{H}^t, \mathbf{X}) \quad (3.23)$$

is used with \mathbf{H}^t denoting iteration t of \mathbf{H} . The computations of f and g can be understood as the feedforward neural network.

To learn the parameters of this GNN, with target information t_v for a specific node v , the loss is defined as

$$loss = \sum_{i=1}^p (t_i - \mathbf{o}_i) \quad (3.24)$$

where p are the supervised nodes. A gradient-descent strategy is utilized in the learning algorithm, which consist of the following three steps: the states h_v^t are updated iteratively using equation (3.19) until time step T . We then obtain an approximate fixed point solution of equation (3.21): $\mathbf{H}(T) \approx \mathbf{H}$. For the next step the gradients of the weights W are calculated from the loss. Finally, the weights W are updated according to the computed gradient. This allows us to train a model for specific supervised or semi-supervised tasks, referred to as downstream task, and get hidden states of nodes in a graph.

A common practice is stacking k GNN layers so that each node is able to aggregate information from nodes within its k -hop neighborhood, seeking an increase in performance. It is important to note that this approach may also increase the noisy information spread by the exponentially increasing neighborhood nodes[16]. In the following, we briefly present two models that are used in the course of this work.

TODO: include figure of graph with neighborhood?

Graph Convolutional Network

Graph convolutional networks (GCNs) aim to generalize the convolution operation of convolutional neural networks (CNNs) to the graph domain. An example is the model proposed by Kipf et al.[14] that introduces a simple, layer-wise propagation rule for multi-layer GCNs as

$$H^{l+1} = \sigma(\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} H^{(l)} W^{(l)}), \quad (3.25)$$

where $\tilde{A} = A + I_N$ is the adjacency matrix of the undirected input graph G with added self-connections. I_N is the identity matrix, \tilde{D} is the degree matrix and $W^{(l)}$ is a layer-specific trainable weight matrix. $\sigma(\cdot)$ denotes an activation function, such as ReLU. $H^{(l)} \in \mathbb{R}^{n \times d}$ denotes the matrix of node activations in the l -th layer, where $H^{(0)} = X$. Note that this definition differs slightly from the iterative formulation in equation 3.23, in the sense that each layer applies a different function F_l , rather than a shared function F . $\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}}$ describes a normalization of the adjacency matrix.

TODO: generalization to signal X ?

$$Z = \hat{D}^{-\frac{1}{2}} \hat{A} \hat{D}^{-\frac{1}{2}} X \Theta \quad (3.26)$$

Higher-order GNN

TODO: A basic GNN model can be implemented as follows (Hamilton, Ying, and Leskovec 2017b) Morris et al.[20] propose an implementation for a GNN model consisting of stacked neural network layers, that each aggregate the local neighborhood information of a node and pass it to the next one. This network is defined specifically for graphs that can be partitioned into r color classes and therefore applicable to bipartite graphs. The new features of node i are then computed with:

$$x_i^{(l)} = \sigma(W_1^{(l)} x_i^{(l-1)} + W_2^{(l)} \cdot \sum_{j \in \mathcal{N}(i)} e_{j,i} \cdot x_j^{(l-1)}), \quad (3.27)$$

where $W_1^{(l)}$ and $W_2^{(l)}$ are two layer-specific trainable weight matrices.

3.5 Perturbation-based Explainability in GNNs

Methods in DL have seen growth in performance in many tasks of artificial intelligence, including GNNs. However, the interpretability of these models is often limited due to their black-box design. Explainability methods aim to bypass this limitation by designing post-hoc techniques that provide insights into the decision-making process in the form of explanations. Such human-intelligible explanations are crucial for deploying models in real-world applications, especially when applied in interdisciplinary fields. There exist several different approaches for explaining predictions of deep graph models, that can be categorized into instance-level methods and model-level methods (see Yuan et al. [29]). Instance-level methods aim to explain each input-graph by identifying important input features for its prediction, leading to input-dependent explanations. These can further be grouped by their importance score calculation into gradients/feature-based, perturbation-based, decomposition methods and surrogate methods. Model-level methods, on the other hand, aim to explain GNNs without considering specific inputs, leading to input-independent, high-level explanations.

In this work we focus on the perturbation-based approach, more specifically the PGExplainer[17], that aims to evaluate the change of prediction with respect to input perturbations. The intuition behind this is that when input information crucial to the prediction is kept, the new prediction should roughly align with the prediction from the original input. The general pipeline for different perturbation based approaches can be described as follows: First, the important features from the input graph are converted into a mask by our generation algorithm, depending on the explanation task at hand. These masks are applied to the input graph to highlight said features. Lastly, the masked graph is fed into the trained GNN to evaluate the mask and update the mask generation algorithm according to the similarity of the predictions. These different approaches mostly differ in the specific mask generation algorithm, the type of mask used and the objective function. It is important to distinguish between soft masks, discrete masks and approximated dis-

crete masks. Soft masks take continuous values between $[0, 1]$ which enables the graph algorithm to be updated via backpropagation. A downside of soft masks is that they suffer from the "introduced evidence" problem (see Dabkowski et al.[6]). Any mask value that is non-zero or non-one may add new semantic meaning or noise to the input graph, since graph edges are by nature discrete. Discrete masks however always rely on non-differentiable operations, e.g. sampling. Thus, the approximated discrete masks utilize reparameterization tricks to avoid the "introduced evidence" problem while also enabling back-propagation.

Explanations can on the one hand be evaluated by visualizing the graph and considering the "human-comprehensibility". Since this requires a ground truth, is prone to the subjective understanding and is usually performed for a few random samples, it is important to apply stable evaluation metrics. One relevant accuracy metric for synthetic datasets with ground truths is the Area Under the Receiver Operating Characteristic Curve (ROC-AUC) (see Richardson et al.[22]). The Receiver Operating Characteristic (ROC) curve plots the False Positive Rate (FPR) on the x-axis against the True Positive Rate (TPR), across different classification thresholds. The area under the curve (AUC) is calculated for said curve, resulting in the ROC-AUC. It is important to note, that a value of 0.5 equals random guessing, while a score of 1.0 indicates perfect classification. TODO: Other metrics include fidelity, results of taxonomy propose only using PGExplainer for Node Classification as it achieves low fidelity on Graph tasks

To fully trust the explanations provided by an explainer model, they must satisfy certain criteria, since there often is a mismatch between the optimizable metrics like accuracy and the actual metric of interest, which may not be measurable (see Ribeiro et al.[21]). First and foremost, an explanation should be **interpretable** and therefore provide qualitative, human-understandable interpretations, that also consider the possibility of limited user knowledge. Additionally, **local fidelity** asserts that explanations should be faithful in a local context and consider the models' behavior in the vicinity of predicted instances. Explainers that treat the model to be explained as a black-box are **model-agnostic** and should therefore be able to explain any model. Lastly, a **global perspective** is needed to explain a model fully, allowing us to take sample explanations of individual predictions that serve as representation of the model.

TODO: figure of perturbation pipeline?

3.6 Boolean Satisfiability Problem

We define the Boolean Satisfiability Problem (SAT) according to Guo et al.[11][p.641]: A Boolean formula is constructed from Boolean variables, that only evaluate to True (1) or False (0), and the three logic operators conjunction (\wedge), disjunction (\vee) and negation (\neg). SAT aims to evaluate whether there exists a variable assignment for a formula constructed of said parts so that it evaluates to True. If so, the formula is said to be satisfiable or

unsatisfiable otherwise. Every propositional formula can be converted into an equivalent formula in conjunctive normal form (CNF), which consists of a conjunction of one or more clauses. These clauses must contain only disjunctions of at least one literal (a variable or its negation). In this work we consider only formulas in CNF, as NeuroSAT[25] assumes SAT problems to be in CNF. An example of a satisfiable formula in CNF over the set of variables $V = \{x_1, x_2\}$ is

$$\psi(V) = (x_1) \wedge (\neg x_1 \vee x_2) \wedge (\neg x_2 \vee x_2)$$

with satisfying assignment $A : \{x_1 \mapsto 1, x_2 \mapsto 1\}$. Furthermore, SAT is *NP*-complete, meaning that if there exists a deterministic algorithm able to solve SAT in polynomial time, then such an algorithm exists for every *NP* problem (see Cook[4]). Current state-of-the-art SAT solvers apply searching based methods such as Conflict Driven Clause Learning[19] or Stochastic Local Search[24] with exponential worst-case complexity.

3.6.1 Representation as Bipartite Graph

SAT has extensively been studied in the form of graphs. Guo et al.[11] describe four different types of graph representations for CNF formulae with varying complexity and information compression. Since we want to minimize the loss of information for SAT we adapt the information-richest form of a literal-clause graph (LCG). A LCG is a bipartite graph that separates literals and clauses, with edges connecting literals to the clauses they appear in. The resulting graph can formally be described by a biadjacency matrix B of shape $l \times c$.

Let $A \in \mathbb{R}^{l+c \times l+c}$ be the adjacency matrix of our bipartite graph. Since for the bipartite case edges exist only between the two color classes l and c , the adjacency matrix can be represented as

$$A(i, j) = \begin{bmatrix} 0_{l \times l} & B \\ B^T & 0_{c \times c} \end{bmatrix}, \quad (3.28)$$

where 0 denotes a zero matrix in the shape of their subscript (see Sun et al.[26]).

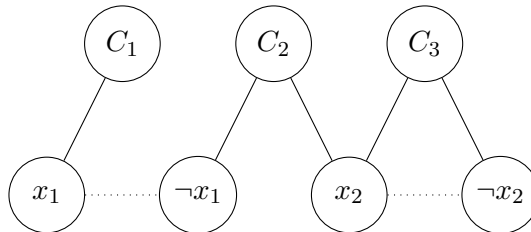


Figure 3.2: LCG representation of $\psi(V)$ with dashed lines representing the connection between complementary literals relevant for the message passing in GNNs.

3.6.2 Unsatisfiable Cores

The core of an unsatisfiable formula in CNF is a subset of the formula that is also unsatisfiable. Every unsatisfiable formula therefore is a core on its own, but can be broken down into smaller cores. The smaller a core the more significance it holds. A minimal unsatisfiable core is also referred to as a minimal unsatisfiable subset (MUS). SAT solvers like minisat[8] are able to compute unsatisfiable cores but do not generally provide a MUS due to high computational cost. However, several deletion-based algorithms exist for computing MUSs (see Torlak et al.[27]).

Chapter 4

PGExplainer - Main part

TODO: This can be understood as an extension of the theoretical background?

V1: In the following chapter, we introduce the PGExplainer[17] and its concepts. The idea is to generate explanations in the form of probabilistic graph generative models for any learned GNN model, henceforth referred to as the downstream task (DT), by utilizing a deep neural network to parameterize the generation process. This approach seeks to explain multiple instances collectively, as they share the same neural network parameters, and therefore improve the generalizability to previous works, particularly the GNNEExplainer[28].

V2: In the following chapter, we introduce the PGExplainer[17] and its concepts. The idea is to generate explanations in the form of probabilistic graph generative models that have proven to learn the concise underlying structures of GNNs most relevant to their predictions. This approach may be applied to any learned GNN model, henceforth referred to as the downstream task (DT), by utilizing a deep neural network to parameterize the generation process. PGExplainer seeks to explain multiple instances collectively, as they share the same neural network parameters, and therefore improve the generalizability to previous works, particularly the GNNEExplainer[28]. This means that all edges in the dataset are predicted by the same model, which leads to a global understanding of the DT.

V3: In the following chapter, we introduce the PGExplainer by Luo et al.[17] and its concepts. The idea is to generate explanations in the form of edge distributions or soft masks using a probabilistic generative model for graph data, known for being able to learn the concise underlying structures from the observed graph data. The explainer uncovers said underlying structures, believed to have the biggest impact on the prediction of a GNNs, as explanations. This approach may be applied to any trained GNN model, henceforth referred to as the target model (TM). By utilizing a deep neural network to parameterize the generation process, the explainer learns to collectively explain multiple instances of a model. Since the parameters of the neural network are shared across the population of explained instance, PGExplainer provides "model-level explanations for each

instance with a global view of the GNN model”. Furthermore, this approach cannot only be used in a transductive setting, but also in an inductive setting, where explanations for unexplained nodes can be generated without retraining the explanation model. This improves the generalizability compared to previous works, particularly the GNNExplainer by Ying et al.[28].

The focus in this approach lies in explaining the graph structure, rather than the graph features, as feature explanations are already common in non-graph neural networks.

Transductive/Inductive explanation?

We then describe our reimplementation in detail (Section 3.2), including the changes made and difficulties during the process.

In Section 3.3 we present the idea of applying PGExplainer on the NeuroSAT framework to generate explanations for the machine learning SAT-solving approach and comparing these to ”human-understandable” concepts like UNSAT cores and backbone variables.

4.1 Theory

We follow the structure of the original paper[[<empty citation>](#)] and start by describing the learning objective (Section 3.1.1), the reparameterization trick (Section 3.1.2), the global explanations (Section 3.1.3) and regularization terms (Section 3.1.4).

4.1.1 Learning Objective

To explain the predictions made by a GNN model for an original input graph G_o with m edges we first define the graph as a combination of two subgraphs: $G_o = G_s + \Delta G$, where G_s represents the subgraph holding the most relevant information for the prediction of a GNN, referred to as explanatory graph. ΔG contains the remaining edges that are deemed irrelevant for the prediction of the GNN. Inspired by GNNExplainer[28], the PGExplainer then finds G_s by maximizing the mutual information between the predictions of the target model and the underlying G_s :

$$\max_{G_s} I(Y_o; G_s) = H(Y_o) - H(Y_o | G = G_s), \quad (4.1)$$

where Y_o is the prediction of the target model with G_o as input. This quantifies the probability of prediction Y_o when the input graph is restricted to the explanatory graph G_s , as in the case of $I(Y_o; G_s) = 1$, knowing the explanatory graph G_s gives us complete information about Y_o , and vice versa. Intuitively, if removing an edge (i, j) changes the prediction of a GNN drastically, this edge is considered important and should therefore be included in G_s . This idea originates from traditional forward propagation based methods

for whitebox explanations (see Dabkowski et al.[6]). It is important to note that $H(Y_o)$ is only related to the target model with fixed parameters during the evaluation/explanation stage. This leads to the objective being equivalent to minimizing the conditional entropy $H(Y_o|G = G_s)$.

To optimize this function a relaxation is applied for the edges, since normally there would be 2^m candidates for G_s . The explanatory graph is henceforth assumed to be a Gilbert random graph, where the selections of edges from G_o are conditionally independent to each other. However, the authors describe a random graph with each edge having its own probability, rather than a shared probability as described in 3.2, as follows: Let $e_{ij} \in V \times V$ be the binary variable indicating whether the edge is selected, with $e_{ij} = 1$ if edge (i, j) is selected to be in the graph, and 0 otherwise. For the random graph variable G the probability of a graph G can be factorized as

$$P(G) = \prod_{(i,j) \in E} P(e_{ij}). \quad (4.2)$$

TODO: Inhomogeneous Erdos Renyi model? Mention that this is a generative model? (A Gilbert random graph is an example of a generative probabilistic model on graph data?) $P(e_{ij})$ is instantiated with the Bernoulli distribution $e_{ij} \sim \text{Bern}(\theta_{ij})$, where $P(e_{ij} = 1) = \theta_{ij}$ is the probability that edge (i, j) exists in G . After this relaxation the learning objective becomes:

$$\min_{G_s} H(Y_o|G = G_s) = \min_{G_s} \mathbb{E}_{G_s} [H(Y_o|G = G_s)] \approx \min_{\Theta} \mathbb{E}_{G_s \sim q(\Theta)} [H(Y_o|G = G_s)], \quad (4.3)$$

where $q(\Theta)$ is the distribution of the explanatory graph that is parameterized by Θ 's.

4.1.2 Reparameterization Trick

As described in section 3.5, a reparameterization trick can be utilized to relax discrete edge weights to continuous variables in the range $(0, 1)$. PGExplainer uses the reparameterizable Gumbel-Softmax estimator[13] to allow for efficiently optimizing the objective function with gradient-based methods. This method introduces the Gumbel-Softmax distribution, a continuous distribution used to approximate samples from a categorical distribution. A temperature τ is used to control the approximation, usually starting from a high value and annealing to a small, non-zero value. Samples with $\tau > 0$ are not identical to samples from the corresponding continuous distribution, but are differentiable and therefore allow back-propagation. The sampling process $G_s \sim q(\Theta)$ of PGExplainer is therefore approximated with a determinant function that takes as input the parameters Ω , a temperature τ and an independent random variable ϵ : $G_s \approx \hat{G} = f_{\Omega}(G_o, \tau, \epsilon)$. The binary concrete distribution[18], also referred to as Gumbel-Softmax distribution, is utilized as an instantiation for the sampling, yielding the weight $\hat{e}_{ij} \in (0, 1)$ for edge (i, j) in \hat{G}_s , computed by:

$$\epsilon \sim \text{Uniform}(0, 1), \quad \hat{e}_{ij} = \sigma((\log \epsilon - \log(1 - \epsilon) + \omega_{ij})/\tau), \quad (4.4)$$

where $\sigma(\cdot)$ is the Sigmoid function and $\omega_{ij} \in \mathbb{R}$ is an explainer logit for the corresponding edge used as a parameter. When $\tau \rightarrow 0$, e.g. during the explanation stage, the weight \hat{e}_{ij} is binarized with the sigmoid function $\lim_{\tau \rightarrow 0} P(\hat{e}_{ij} = 1) = \frac{\exp(\omega_{ij})}{1 + \exp(\omega_{ij})}$. Since $P(e_{ij} = 1) = \theta_{ij}$, choosing $\omega_{ij} = \log \frac{\theta_{ij}}{1 - \theta_{ij}}$ leads to $\lim_{\tau \rightarrow 0} \hat{G}_s = G_s$ and justifies the approximation of the Bernoulli distribution with the binary concrete distribution. During training, when $\tau > 0$, the objective function in (4.3) is smoothed with a well-defined gradient $\frac{\partial \hat{e}_{ij}}{\partial \omega_{ij}}$ and becomes:

$$\min_{\Omega} \mathbb{E}_{\epsilon \sim \text{Uniform}(0,1)} H(Y_o | G = \hat{G}_s) \quad (4.5)$$

The authors follow the approach of GNNExplainer[28] and modify the objective by replacing the conditional entropy with cross entropy between the label class and the prediction of the target model. This is justified by the greater importance of understanding the model’s prediction of a certain class, rather than providing an explanation based solely on its confidence.

With the modification to cross-entropy $H(Y_o, \hat{Y}_s)$, where \hat{Y}_s is the prediction of the target model when \hat{G}_s is given as input, as well as the adaption of Monte Carlo sampling, the learning objective becomes: TODO: ONE LABEL FOR BOTH EQUATIONS!?

$$\min_{\Omega} \mathbb{E}_{\epsilon \sim \text{Uniform}(0,1)} H(Y_o, \hat{Y}_s) \approx \min_{\Omega} -\frac{1}{K} \sum_{k=1}^K \sum_{c=1}^C P(Y_o = c) \log P(\hat{Y}_s = c) \quad (4.6)$$

$$= \min_{\Omega} -\frac{1}{K} \sum_{k=1}^K \sum_{c=1}^C P_{\Phi}(Y_o = c | G = G_o) \log P_{\Phi}(\hat{Y}_s = c | G = \hat{G}_s^{(k)}). \quad (4.7)$$

Φ denotes the parameters in the target model, K is the number of total sampled graphs, C is the number of class labels, and $\hat{G}_s^{(k)}$ denotes the k -th graph sampled with equation 4.4, parameterized by Ω .

4.1.3 Global Explanations

The novelty of PGExplainer lies in the ability to generate explanations for graph data with a global perspective, that allow for understanding the general picture of a model across a population. This saves resources when analyzing large graph datasets, as new instances can be explained without retraining the model, and can also be helpful for establishing the users’ trust in these explanations. To achieve this the authors propose the use of a parameterized network that learns to generate explanations from the target model, which also apply to not yet explained instances. PGExplainer hence explains predictions of the target model over multiple instances collectively (TODO: Cut this sentence?).

Since GNNs apply two functions F and G to calculate the global state embeddings and downstream task outputs respectively, we denote these two functions as $\text{GNNE}_{\Phi_0}(\cdot)$ and $\text{GNNC}_{\Phi_1}(\cdot)$ for any GNN in the context of PGExplainer. For models without explicit classification layers the last layer is used to compute the output instead. It follows

$$\mathbf{Z} = \text{GNNE}_{\Phi_0}(G_o, \mathbf{X}), \quad Y = \text{GNNC}_{\Phi_1}(\mathbf{Z}), \quad (4.8)$$

where \mathbf{Z} denotes the matrix of final node representations z , referred to as node embeddings, and the initial state is G_o . TODO: (Because of focus on graph structure rather than features?) For generalizability across different GNN layers the output is only dependent on the node representation, that encapsulates both features and structure of the input graph. This representation also serves as the input for the explainer network g , defined as: TODO: Rename g ?

$$\Omega = g_{\Psi}(G_o, \mathbf{Z}). \quad (4.9)$$

Ψ denotes the parameters in the explanation network and the output Ω is treated as parameter in equation 4.6. Since Ψ is shared by all edges among the population, PGExplainer collectively provides explanations for multiple instances. Thus, the learning objective in a collective setting with \mathcal{I} being the set of instances becomes:

$$\min_{\Psi} -\frac{1}{K} \sum_{i \in \mathcal{I}} \sum_{k=1}^K \sum_{c=1}^C P_{\Phi}(Y_o = c | G = G_o^{(i)}) \log P_{\Phi}(\hat{Y}_s = c | G = \hat{G}_s^{(i,k)}). \quad (4.10)$$

Consequently, $G^{(i)}$ and $G_s^{(i,k)}$ denote the input graph and the k -th graph sampled with equation 4.4 in 4.9 respectively for instance i . The authors propose two slightly different instantiations for node classification and graph classification tasks.

Explanation network for node classification

Let an edge (i, j) be considered relevant for the prediction of a node u , but irrelevant for the prediction of a node v . To explain the prediction of node v we specify the network in 4.9 as:

$$\omega_{ij} = \text{MLP}_{\Psi}([\mathbf{z}_i \oplus \mathbf{z}_j \oplus \mathbf{z}_v]). \quad (4.11)$$

MLP_{Ψ} is an MLP (TODO see ?? for implementation details) parameterized with Ψ and \oplus denotes the concatenation operation. Thus, a concatenation of the node embeddings of nodes i, j and v respectively is fed through the network. The output ω_{ij} is therefore an edge logit, which serves as a parameter in the sampling process.

Note that in their codebase the authors use a concatenation of all hidden representations instead of final node embeddings for node level tasks. For a target GNN consisting of multiple graph layers with

$$\begin{aligned} \mathbf{H}_1 &= F_1(G_o, \mathbf{X}), \\ \mathbf{H}_2 &= F_2(\mathbf{H}_1, \mathbf{X}), \\ &\vdots \\ \mathbf{H}_L &= F_L(\mathbf{H}_{L-1}, \mathbf{X}), \end{aligned}$$

this leads to Z being the matrix of node embeddings z that are computed as:

$$z_i = h_{1,i} \oplus h_{2,i} \oplus \dots \oplus h_{L,i} \quad (4.12)$$

Explanation network for graph classification For graph level tasks the authors consider each graph to be an instance, regardless of specific nodes. The specification for the network thus becomes:

$$\omega_{ij} = \text{MLP}_\Psi([\mathbf{z}_i \oplus \mathbf{z}_j]), \quad (4.13)$$

where for each edge in G_o a concatenation of both its nodes is fed through the MLP.

TODO: Include computational complexity of PGE in comparison to GNN?

This leads to an improved computational complexity when compared to their baseline GNNExplainer, since for one the number of parameters in the explainer does no longer depend on the size of the input graph and since the explainer does not have to be retrained for every unexplained instance.

TODO: Include algorithms?

4.1.4 Regularization Terms

To enhance the preservation of desired properties of explanations the authors propose various regularization terms. These are added to the learning objective, depending on the specific downstream task at hand.

Size and entropy constraints

Inspired by GNNExplainer[28], to obtain compact and precise explanations, a constraint on the size of the explanations is added in the form of $\|\Omega\|_1$, the l_1 norm on latent variables Ω . Additionally, to encourage the discreteness of edge weights, element-wise entropy is added as a constraint:

$$H_{\hat{G}_s} = -\frac{1}{|\varepsilon|} \sum_{(i,j) \in \varepsilon} (\hat{e}_{ij} \log \hat{e}_{ij} + (1 - \hat{e}_{ij}) \log(1 - \hat{e}_{ij})), \quad (4.14)$$

for one explanatory graph \hat{G}_s with ε edges. For the collective setting, this is added as a mean over all instances in \mathcal{I} .

Note that the following two constraints are not used in the original experimental setup, but serve as inspiration for constraints introduced in our NeuroSAT application?? and are therefore included.

Budget constraint

The authors propose the modification of the size constraint to a budget constraint, for a predefined available budget B . Let $|\hat{G}_s| \leq B$, then the budget regularization is defined as:

$$R_b = \text{ReLU}\left(\sum_{(i,j) \in \varepsilon} \hat{e}_{ij} - B\right). \quad (4.15)$$

Note that $R_b = 0$ when the explanatory graph is smaller than the budget. When out of budget, the regularization is similar to that of the size constraint.

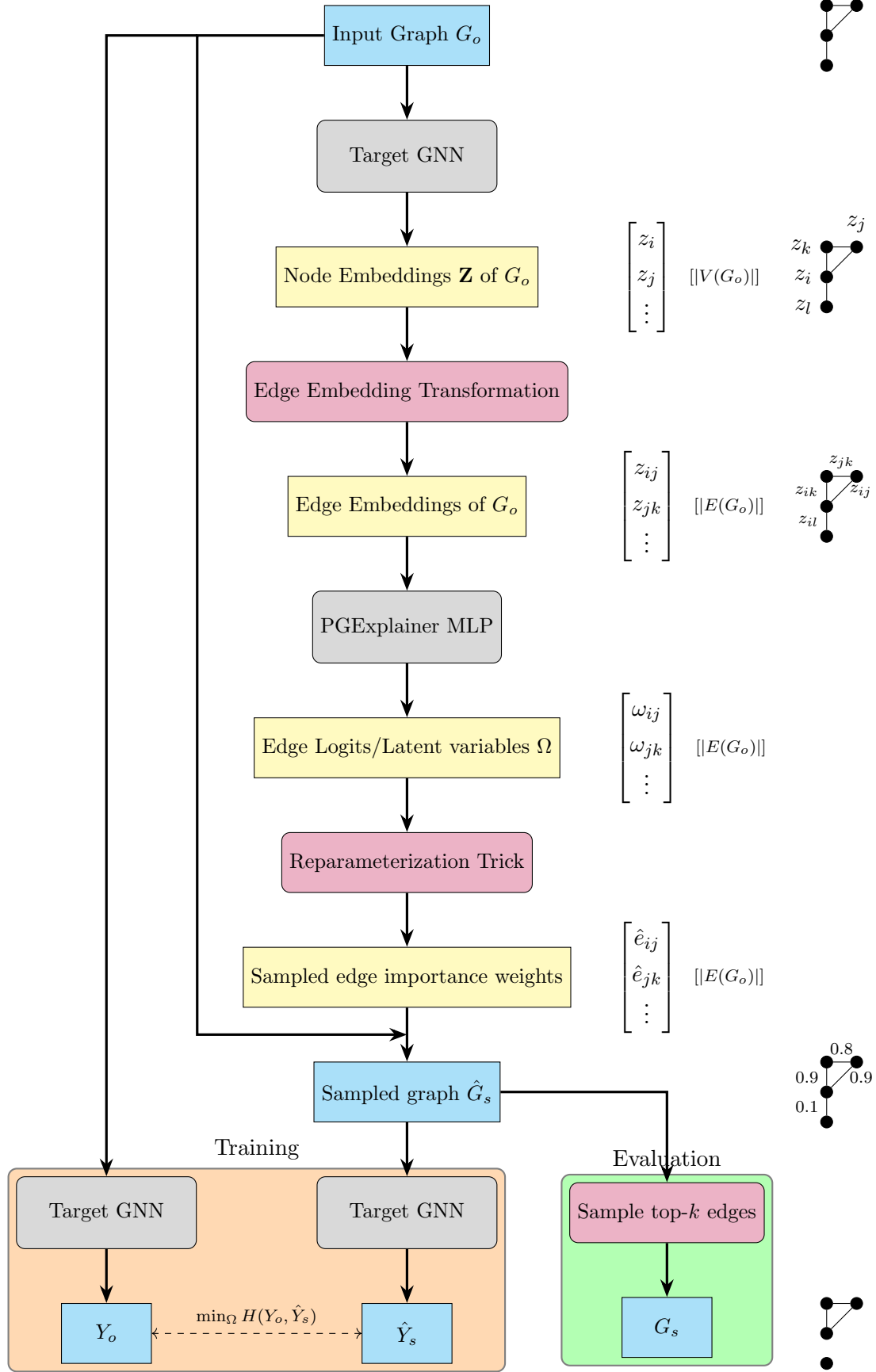


Figure 4.1: The complete pipeline of PGExplainer.

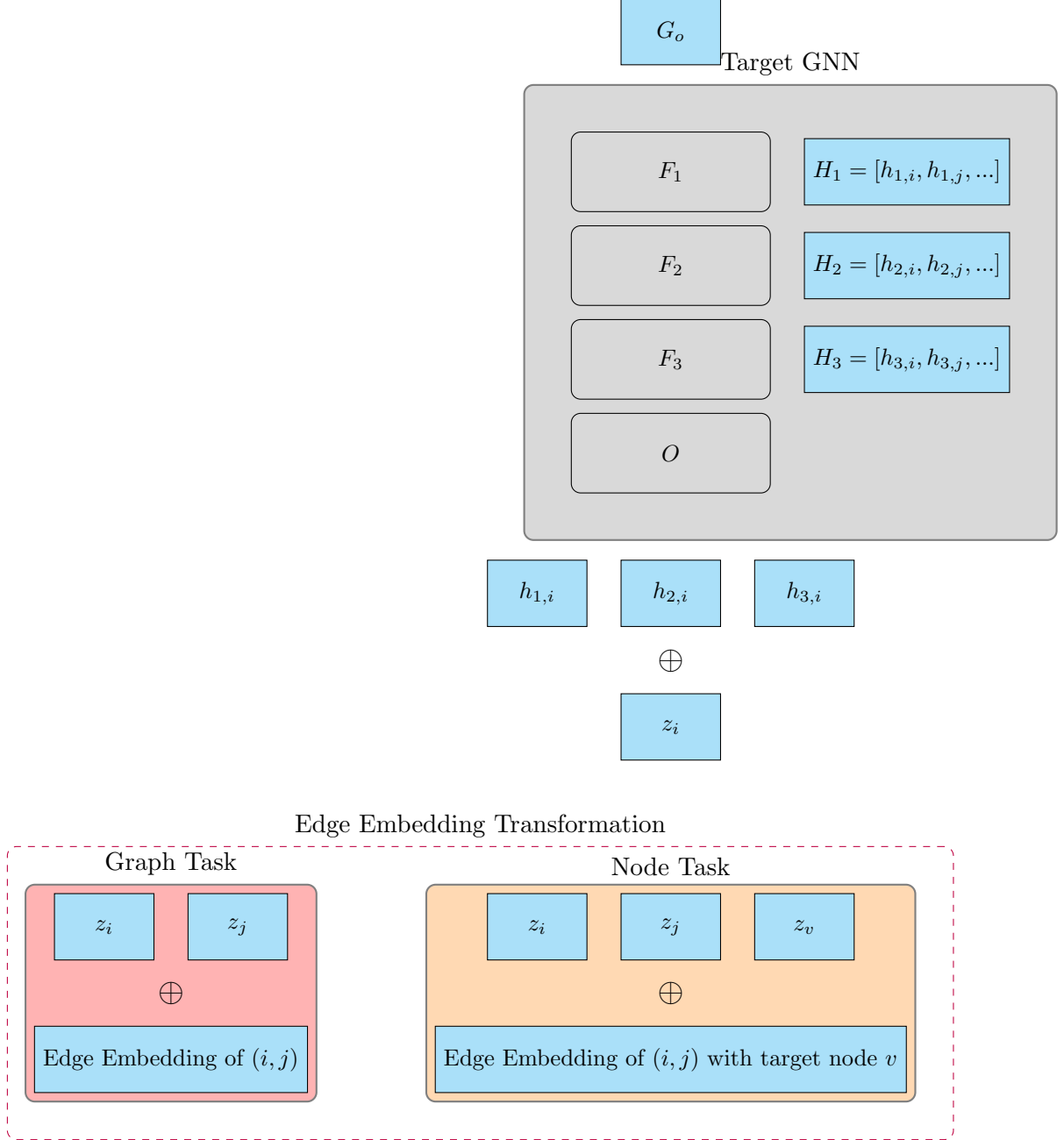


Figure 4.2: TODO: Visualization of the edge embedding transformation used to create inputs for the explainer network. Depending on the downstream task used in the target model the created edge embedding differs slightly.

Connectivity constraint

To enhance the effect of the explainer detecting a small, connected subgraph, motivated through real-life motifs being inherently connected, the authors suggest adding the cross-entropy of adjacent edges. Let (i, j) and (i, k) be two edges that both connect to the node i , then (i, k) should rather be included in the explanatory graph if the edge (i, j) is selected to be included. This is formally defined as:

$$H(\hat{e}_{ij}, \hat{e}_{ik}) = -[1 - \hat{e}_{ij} \log(1 - \hat{e}_{ik}) + \hat{e}_{ij} \log \hat{e}_{ik}]. \quad (4.16)$$

We note that in practice this is implemented only for the two highest edge weights for each edge. The definition therefore would change to (i, j) and (i, k) being the top two edges that connect to node i .

4.2 Reimplementation

Notes on Original code of target model: Not clearly specified in paper, but according to code: - No specific args set for layers etc., except hyperparameters for different datasets (see sweeps) - Therefore, assuming default settings for layers - Graph conv layer uses ReLu activation (done), bias initialized with zeroes (done), no dropout (partially done), and embedding normalization! (TODO), default order 'AW', weight init 'glorot' (Done with Xavier) - WE USED PyG GraphConv since it allows for bipartite+edge weights; but compare to layer close to original!!

Concrete graph layer used in original code is very similar to the semi-supervised layer??:

$$Z = \hat{A}XW \quad (4.17)$$

Note that a pytorch implementation of the semi-supervised layer is used in the replication paper:

$$X' = \hat{D}^{-\frac{1}{2}} \hat{A} \hat{D}^{-\frac{1}{2}} X \Theta \quad (4.18)$$

We use a pytorch implementation of the Weisfeiler and Leman Go Neural layer:

$$x'_i = W_1 x_i + W_2 \sum_{j \in \mathcal{N}(i)} e_{j,i} \cdot x_j \quad (4.19)$$

TODO: Maybe separate the theory of PGE from our own work more strictly? E.g. 3. PGE theory, 4. PGE reimplementation and NeuroSAT application?

Implementation details:

- Started by reimplementing the downstream tasks used in og paper for node and graph class.
- Node datasets taken from original and transformed to a "pyg format", to keep original structure and ground truths

- Graph: BA-2Motif from pyg library with self generated gt, MUTAG dataset taken from pyg and added gt-labels that were added in original dataset
- Created pytorch/pytorch geometric implementation of 3-layer graph conv networks
- architecture as described in paper: ReLu activation, Pooling for graph net
- Xavier initialization for all layers
- Same hyperparameters?(Adam, $1 * 10^{-3}$ lr, 1000 epochs)/experimental setup?
- ADDED Dropout(0.1) to improve performance/overfitting on node tasks
- Fully connected layer: torch geometric GraphConv to allow passing of edge weights(+ bipartite)!?
- Original references sageConv in paper, pyG impl. does not allow edge weights, verify!
- GATConv for Graph attention(referenced by original)
- Alternatives: GCNConv(edge weights, no bipartite graphs)
- 80/10/10 split
- =, Similar accuracies achieved

Explainer: This is irrelevant for paper, description of why reparameterization trick necessary - First approach tried passing a masked edge index to downstream task with edge weights $\in [0, 1]$ - Unable to learn with hard cut-off (bad for gradients). (Same with TopK probably?!) - Pass calculated edge weights to downstream task for learning. If no edge weights are passed(downstream task), all edge weights are initialized with one to represent all edges being relevant

4.3 Application on NeuroSAT

NeuroSAT: Messages are passed between clauses and literals, as well as literals and their complement. 1. Clause receives from neighboring literals 2. Literals receive from clauses and complement.

(Define flip function that swaps literal row with row of its negation; relevant for NeuroSAT) We create required data with provided methods, add unsat cores and MUSs as gt and adapt the NeuroSAT model to allow passing edge weights into the adjacency matrix. PGExplainer adapted for SAT data, NeuroSAT embeddings and evaluation.

Reimplementation of NeuroSAT provided by Rodhi. As the code used for NeuroSAT can also be found in our Repository, we stress that only the changes described in the following chapter are part of our work.

What did we do? What did we change for NeuroSAT? What data was used? How did we adapt PGExplainer?

Only change in NeuroSAT: pass edge weights into adjacency matrix. Calculates
Generated batches of unsat problems that "turned" unsat because of last added clause.
10 literals per problem. Only unsat to test for unsat cores, that only apply for unsat problems. Calculated unsat cores with solver xy by adding negative assumption literals per clause and passing these as assumption for calculation. The edges of the clauses present in the unsat core were treated as ground truth.

Changes for explainer: Edge embeddings calced by DS and passed to explainer. Calculated edge embeddings by concatenating node embeddings for connected nodes, similar to original. Embeddings fed into MLP, weights sampled with reparam. trick to get edge "probabilities", passed as "unbatched" sampled graph into NeuroSAT predictor. Visualization of SAT problems with edge weights, ands gts.

For quant. eval. adapted roc auc as metric as done in PGExplainer. Results seem "good" but qual. eval. shows different result. roc auc bad metric?

For qual. eval. topk(=number of edges in gt) edges of predictions were highlighted to be compared to gt edges. For quant. eval. the edge probabilities were compared to gt with 1s for edges in gt and 0s for rest.

Soft modified connectivity constraint

To account for the definition of UNSAT cores, that consist of sub-clauses of the original combination of clauses in the problem, we add a constraint that reinforces the prediction of complete clauses. Therefore, if the explainer assigns a high score to an edge (i, j) , all edges (j, k) that also connect to the clause node j should receive a high score. Therefore, we introduce a soft constraint that punishes varying edge weights for the same clause. For our sampled bipartite Graph \hat{G}_s with node sets L and C containing literal nodes and clause nodes respectively, we define:

$$R_c = \sum_{c \in C} \text{Var}(E_c) = \sum_{c \in C} \frac{1}{|E(c)|} \sum_{e \in E(c)} (e - \bar{E}_c)^2$$

where $E_c = \{\hat{e}_{i,j} \mid j = c\}$ is the set of edge weights of edges that connect to clause c and \bar{E}_c denotes the mean of E_c .

Hard constraint

Define PGE formulas over grouped clause edges rather than all edges.

The reparameterization trick is still applied for each edge, but ϵ_c is sampled per clause instead of per edge, so that all edges that connect to a clause are forced to not only bear the same logit, but also the exact same weight during training.

$$\epsilon_c \sim \text{Uniform}(0, 1), \quad \hat{e}_{l,c} = \sigma((\log \epsilon_c - \log(1 - \epsilon_c) + \omega_{ij}/\tau)), \quad (4.20)$$

Additionally, we restrain the edge logits $\omega_{i,j}$ calculated by the MLP to be identical for all edges that connect to the same clause. Let $\Omega_c = \{\hat{\omega}_{i,c} \mid (i, c) \in E\}$ denote the set of logits corresponding to the edges connected to node c . We update these logits with:

$$\mu_c = \frac{1}{|\Omega_c|} \sum_{\hat{\omega}_{i,c} \in \Omega_c} \hat{\omega}_{i,c} \quad (4.21)$$

We calculate the mean logit μ_c for all edges that connect to node c with

$$\mu_c = \frac{1}{|E_c|} \sum_{(i,j) \in E_c} \omega_{i,j}. \quad (4.22)$$

The update rule is then defined as

$$\omega'_{i,j} = \mu_c, \quad \forall (i, j) \in E_c \quad (4.23)$$

Chapter 5

Results

Experimental Setup: We follow the experimental setup from the PGExplainer as closely as possible. Since the textual description refers to the setup from GNNExplainer and is lacking in some aspects, we extract the missing information from the codebase. As the hyperparameters are unclear or not comprehensible for some tasks we also draw information from the configs of PYTORCH REIMPL.

We use normalization in our downstream models, though it is not described in the paper, since it is used in the code. We also experimented with the effects of the use of normalization since it seems to be relevant to the performance of the explainer.

The explainer is trained and evaluated on the same data. We also run experiments with added train/test splits TODO!

Not all data is fed into the explainer: BA-Shapes: BA-Community: Tree-Cycles: "First"/to base graph attached node of each motif in graph is used. Tree-Grid: All Motif nodes are used BA-2Motif: All graphs are used MUTAG: Only the graphs with an available ground truth are used. GT exist for mutagenic graphs that have either chemical groups NH2 or NO2. We later discuss if these selections make sense and run experiments with different data selections.

Quantitative: 10 explainer runs on one downstream model; Calculate ROC-AUC over ALL graphs/nodes in each run; Qualitative: Original uses a threshold; we instead take the topK nodes according to the dataset/motif as an explanation

We also discuss if treating the number of k edges as a hyperparameter dependant on the downstream task makes sense and propose having the network learn it, to improve generalizability and allow the explainer to work on data with varying size.

CPU vs GPU:

It is important to highlight that our code achieved better and way more stable results for BA-2Motif when trained on a gpu instead of cpu.

Ba-Shapes, Tree-Cycles and MUTAG results achieved were identical.

Ba-Community and Tree-Grid achieved very slightly better results on CPU.

Sweeps: (Params ordered by importance)

BA-Shapes: higher size reg -i 0.1; lower entropy reg -i 0.01; lr and tT very low impact but slightly higher -i 0.01 and 5. Note that Loss curve jumps on most runs! (logical-sweep-94 and restful-sweep-92 have clean loss) TRY HIGHER SIZE REG AND LOWER ENTROPY REG

BA-Community: lr 0.0001 too low, not working -i 0.003; lower entropy -i 0.1; higher size -i 0.1; TRY MORE SEEDS, LR, EPOCHS?

Tree-Cycles: high lr -i 0.01 ; lower entropy reg -i 0.1/0.01; higher size reg -i 0.1/0.01 ; lower tT -i 1. TRY WITH MORE SEEDS FOR ENT, SIZE, TEMP? Confirmed higher size reg -i 0.1; lower entropy reg -i 0.01; temp really low impact, tendency higher. TRY 30 EPOCHS???

Tree-Grid: high lr -i 0.01; high size reg -i 1; higher entropy reg? -i 10/1, high tT -i 5. TRY MORE SEEDS FOR ENTROPY REG -i not quite clear, tendency lower; MAYBE EVEN HIGHER LR -i No

BA-2Motif: RUN ON GPU - Not the cause. Cause for better results were features of 1 instead of 0.1! However, good results achieved on BA2-Motif dataset from pyg, not original one.

Comparison to original one: Original dataset transformed to pytorch performs way worse, for features of 0.1! Mean AUC of about 0.4!

Original dataset with features changed to ones instead of 0.1: Works good as well.

MUTAG: Low lr -i 0.0003; low entropy reg(high impact, but highest AUC runs vary) -i 0.1; low tT -i 1; less epochs -i 20; low size reg -i 0.005(/0.01); Loss is messy and AUC seems to decrease over time! lr 0.0001 worse, entropy reg 0.1/0.01 has zero effect -i 0.1

Effects of selected motif nodes for Node task: Compare Tree-Grid/Tree-Cycles performance when using all/one node per motif...

Chapter 6

Discussion

Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua.

Chapter 7

Conclusion

Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua.

Bibliography

- [1] A.S. Asratian, T.M.J. Denley, and R. Häggkvist. *Bipartite Graphs and their Applications*. Cambridge Tracts in Mathematics. Cambridge University Press, 1998. ISBN: 9781316582688. URL: <https://books.google.de/books?id=18fLCgAAQBAJ>.
- [2] Hongyun Cai, Vincent W Zheng, and Kevin Chen-Chuan Chang. “A comprehensive survey of graph embedding: Problems, techniques, and applications”. In: *IEEE transactions on knowledge and data engineering* 30.9 (2018), pp. 1616–1637.
- [3] Augustin Cauchy et al. “Méthode générale pour la résolution des systemes d’équations simultanées”. In: *Comp. Rend. Sci. Paris* 25.1847 (1847), pp. 536–538.
- [4] Stephen A Cook. “The complexity of theorem-proving procedures”. In: *Logic, automata, and computational complexity: The works of Stephen A. Cook*. 2023, pp. 143–152.
- [5] T.M. Cover and J.A. Thomas. “Entropy, Relative Entropy, and Mutual Information”. In: *Elements of Information Theory*. John Wiley & Sons, Ltd, 2005, pp. 13–55. ISBN: 9780471748823. DOI: <https://doi.org/10.1002/047174882X.ch2>. eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/047174882X.ch2>. URL: <https://onlinelibrary.wiley.com/doi/abs/10.1002/047174882X.ch2>.
- [6] Piotr Dabkowski and Yarin Gal. “Real time image saliency for black box classifiers”. In: *Advances in neural information processing systems* 30 (2017).
- [7] Reinhard Diestel. “Random Graphs”. In: *Graph Theory*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2017, pp. 323–345. ISBN: 978-3-662-53622-3. DOI: 10.1007/978-3-662-53622-3_11. URL: https://doi.org/10.1007/978-3-662-53622-3_11.
- [8] Niklas Eén and Niklas Sörensson. “An extensible SAT-solver”. In: *International conference on theory and applications of satisfiability testing*. Springer. 2003, pp. 502–518.
- [9] Edgar N Gilbert. “Random graphs”. In: *The Annals of Mathematical Statistics* 30.4 (1959), pp. 1141–1144.
- [10] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. <http://www.deeplearningbook.org>. MIT Press, 2016.

- [11] Wenxuan Guo et al. “Machine learning methods in solving the boolean satisfiability problem”. In: *Machine Intelligence Research* 20.5 (2023), pp. 640–655.
- [12] Lars Holdijk et al. “[Re] Parameterized Explainer for Graph Neural Network”. In: *ML Reproducibility Challenge 2020*. 2021.
- [13] Eric Jang, Shixiang Gu, and Ben Poole. “Categorical reparameterization with gumbel-softmax”. In: *arXiv preprint arXiv:1611.01144* (2016).
- [14] Thomas N Kipf and Max Welling. “Semi-supervised classification with graph convolutional networks”. In: *arXiv preprint arXiv:1609.02907* (2016).
- [15] Y. Lecun et al. “Gradient-based learning applied to document recognition”. In: *Proceedings of the IEEE* 86.11 (1998), pp. 2278–2324. DOI: 10.1109/5.726791.
- [16] Zhiyuan Liu and Jie Zhou. “Introduction”. In: *Introduction to Graph Neural Networks*. Cham: Springer International Publishing, 2020, pp. 1–3. ISBN: 978-3-031-01587-8. DOI: 10.1007/978-3-031-01587-8_1. URL: https://doi.org/10.1007/978-3-031-01587-8_1.
- [17] Dongsheng Luo et al. “Parameterized explainer for graph neural network”. In: *Advances in neural information processing systems* 33 (2020), pp. 19620–19631.
- [18] Chris J Maddison, Andriy Mnih, and Yee Whye Teh. “The concrete distribution: A continuous relaxation of discrete random variables”. In: *arXiv preprint arXiv:1611.00712* (2016).
- [19] Joao P Marques-Silva and Karem A Sakallah. “GRASP: A search algorithm for propositional satisfiability”. In: *IEEE Transactions on Computers* 48.5 (1999), pp. 506–521.
- [20] Christopher Morris et al. “Weisfeiler and leman go neural: Higher-order graph neural networks”. In: *Proceedings of the AAAI conference on artificial intelligence*. Vol. 33. 01. 2019, pp. 4602–4609.
- [21] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. “” Why should i trust you?” Explaining the predictions of any classifier”. In: *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*. 2016, pp. 1135–1144.
- [22] Eve Richardson et al. “The receiver operating characteristic curve accurately assesses imbalanced datasets”. In: *Patterns* 5.6 (2024), p. 100994. ISSN: 2666-3899. DOI: <https://doi.org/10.1016/j.patter.2024.100994>. URL: <https://www.sciencedirect.com/science/article/pii/S2666389924001090>.
- [23] Franco Scarselli et al. “The Graph Neural Network Model”. In: *IEEE Transactions on Neural Networks* 20.1 (2009), pp. 61–80. DOI: 10.1109/TNN.2008.2005605.
- [24] Bart Selman, Henry A Kautz, Bram Cohen, et al. “Local search strategies for satisfiability testing.” In: *Cliques, coloring, and satisfiability* 26 (1993), pp. 521–532.

- [25] Daniel Selsam et al. “Learning a SAT solver from single-bit supervision”. In: *arXiv preprint arXiv:1802.03685* (2018).
- [26] Hanqing Sun et al. “Supervised biadjacency networks for stereo matching”. In: *Multimedia Tools and Applications* 83 (June 2023), pp. 1–26. DOI: 10.1007/s11042-023-15362-5.
- [27] Emina Torlak, Felix Sheng-Ho Chang, and Daniel Jackson. “Finding Minimal Unsatisfiable Cores of Declarative Specifications”. In: *FM 2008: Formal Methods*. Ed. by Jorge Cuellar, Tom Maibaum, and Kaisa Sere. Berlin, Heidelberg: Springer Berlin Heidelberg, 2008, pp. 326–341. ISBN: 978-3-540-68237-0.
- [28] Zhitao Ying et al. “Gnnexplainer: Generating explanations for graph neural networks”. In: *Advances in neural information processing systems* 32 (2019).
- [29] Hao Yuan et al. “Explainability in graph neural networks: A taxonomic survey”. In: *IEEE transactions on pattern analysis and machine intelligence* 45.5 (2022), pp. 5782–5799.

Appendix A

Supplementary Material

Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua.

Experiment XY

Eidesstattliche Versicherung

(Affidavit)

Name, Vorname
(surname, first name)

Matrikelnummer
(student ID number)

☐ Bachelorarbeit
(Bachelor's thesis)

☐ Masterarbeit
(Master's thesis)

Titel
(Title)

Ich versichere hiermit an Eides statt, dass ich die vorliegende Abschlussarbeit mit dem oben genannten Titel selbstständig und ohne unzulässige fremde Hilfe erbracht habe. Ich habe keine anderen als die angegebenen Quellen und Hilfsmittel benutzt sowie wörtliche und sinngemäße Zitate kenntlich gemacht. Die Arbeit hat in gleicher oder ähnlicher Form noch keiner Prüfungsbehörde vorgelegen.

I declare in lieu of oath that I have completed the present thesis with the above-mentioned title independently and without any unauthorized assistance. I have not used any other sources or aids than the ones listed and have documented quotations and paraphrases as such. The thesis in its current or similar version has not been submitted to an auditing institution before.

Ort, Datum
(place, date)

Unterschrift
(signature)

Belehrung:

Wer vorsätzlich gegen eine die Täuschung über Prüfungsleistungen betreffende Regelung einer Hochschulprüfungsordnung verstößt, handelt ordnungswidrig. Die Ordnungswidrigkeit kann mit einer Geldbuße von bis zu 50.000,00 € geahndet werden. Zuständige Verwaltungsbehörde für die Verfolgung und Ahndung von Ordnungswidrigkeiten ist der Kanzler/die Kanzlerin der Technischen Universität Dortmund. Im Falle eines mehrfachen oder sonstigen schwerwiegenden Täuschungsversuches kann der Prüfling zudem exmatrikuliert werden. (§ 63 Abs. 5 Hochschulgesetz - HG -).

Die Abgabe einer falschen Versicherung an Eides statt wird mit Freiheitsstrafe bis zu 3 Jahren oder mit Geldstrafe bestraft.

Die Technische Universität Dortmund wird ggf. elektronische Vergleichswerkzeuge (wie z.B. die Software „turnitin“) zur Überprüfung von Ordnungswidrigkeiten in Prüfungsverfahren nutzen.

Die oben stehende Belehrung habe ich zur Kenntnis genommen:

Official notification:

Any person who intentionally breaches any regulation of university examination regulations relating to deception in examination performance is acting improperly. This offense can be punished with a fine of up to EUR 50,000.00. The competent administrative authority for the pursuit and prosecution of offenses of this type is the Chancellor of TU Dortmund University. In the case of multiple or other serious attempts at deception, the examinee can also be unenrolled, Section 63 (5) North Rhine-Westphalia Higher Education Act (*Hochschulgesetz, HG*).

The submission of a false affidavit will be punished with a prison sentence of up to three years or a fine.

As may be necessary, TU Dortmund University will make use of electronic plagiarism-prevention tools (e.g. the "turnitin" service) in order to monitor violations during the examination procedures.

I have taken note of the above official notification:*

Ort, Datum
(place, date)

Unterschrift
(signature)

***Please be aware that solely the German version of the affidavit ("Eidesstattliche Versicherung") for the Bachelor's/ Master's thesis is the official and legally binding version.**