

## Explainable Graph Attention Networks

**Authors :** David Pham, Yongfeng Zhang

**Abstract :** Graphs are an important structure for data storage and computation. Recent years have seen the success of deep learning on graphs such as Graph Neural Networks (GNN) on various data mining and machine learning tasks. However, most of the deep learning models on graphs cannot easily explain their predictions and are thus often labelled as “black boxes.” For example, Graph Attention Network (GAT) is a frequently used GNN architecture, which adopts an attention mechanism to carefully select the neighborhood nodes for message passing and aggregation. However, it is difficult to explain why certain neighbors are selected while others are not and how the selected neighbors contribute to the final classification result. In this paper, we present a graph learning model called Explainable Graph Attention Network (XGAT), which integrates graph attention modeling and explainability. We use a single model to target both the accuracy and explainability of problem spaces and show that in the context of graph attention modeling, we can design a unified neighborhood selection strategy that selects appropriate neighbor nodes for both better accuracy and enhanced explainability. To justify this, we conduct extensive experiments to better understand the behavior of our model under different conditions and show an increase in both accuracy and explainability.

**Keywords :** explainable AI, graph attention network, graph neural network, node classification

**Conference Title :** ICPRL 2023 : International Conference on Pattern Recognition and Machine Learning

**Conference Location :** Boston, United States

**Conference Dates :** April 17-18, 2023