# Interpretable Graph Neural Networks for Tabular Data

**Amr Alkhatib**[1], **Sofiane Ennadir**[1], **Henrik Boström**[1], **Michalis Vazirgiannis**[1,2]

[1]KTH Royal Institute of Technology
Electrum 229, 164 40 Kista, Stockholm, Sweden
{alkhat,ennadir,bostromh,mvaz}@kth.se
[2]DaSciM, LIX, École Polytechnique, Institut Polytechnique de Paris, France.
mvazirg@lix.polytechnique.fr

## Abstract

Data in tabular format is frequently occurring in real-world applications. Graph Neural Networks (GNNs) have recently been extended to effectively handle such data, allowing feature interactions to be captured through representation learning. However, these approaches essentially produce black-box models, in the form of deep neural networks, precluding users from following the logic behind the model predictions. We propose an approach, called IGNNet (Interpretable Graph Neural Network for tabular data), which constrains the learning algorithm to produce an interpretable model, where the model shows how the predictions are exactly computed from the original input features. A large-scale empirical investigation is presented, showing that IGNNet is performing on par with state-of-the-art machine-learning algorithms that target tabular data, including XGBoost, Random Forests, and TabNet. At the same time, the results show that the explanations obtained from IGNNet are aligned with the true Shapley values of the features without incurring any additional computational overhead.

## 1    Introduction

In some application domains, e.g., medicine and law, predictions made by machine learning models need justification for legal and ethical considerations (Lakkaraju et al. 2017; Goodman and Flaxman 2017). In addition, users may put trust in such models only with a proper understanding of the reasoning behind the predictions. A direct solution is to use learning algorithms that produce interpretable models, such as logistic regression (Berkson 1944), which provides both local (instance-specific) and global (model-level) explanations for the predictions. However, such algorithms often result in a substantial loss in predictive performance compared to algorithms that generate black-box models, e.g., XGBoost (Chen and Guestrin 2016), Random Forests (Breiman 2001), and deep learning algorithms (Pintelas, Livieris, and Pintelas 2020; Mori and Uchihira 2019). Post-hoc explanation techniques, e.g., SHAP (Lundberg and Lee 2017), LIME (Ribeiro, Singh, and Guestrin 2016), and Anchors (Ribeiro, Singh, and Guestrin 2018), have been put forward as tools to explain predictions of the black-box models. However, the explanations provided by such techniques are limited in that they either do not show how exactly the predictions are computed, but merely present feature scores, such as LIME and SHAP, or come with no

guarantees on the fidelity, i.e., that the provided explanation agrees with the underlying model (Yeh et al. 2019; Delaunay, Galárraga, and Largouët 2020). As extensively argued in (Rudin 2019), there are hence several reasons to consider generating interpretable models in the first place, if trustworthiness is a central concern.

Graph Neural Networks (GNNs) have emerged as a powerful framework for representation learning of graph-structured data (Xu et al. 2019). The application of GNNs has been extended to tabular data, where a GNN can be used to learn an enhanced representation for the data points (rows) or to model the interaction between different features (columns). TabGNN (Guo et al. 2021) is an example of the first approach, where each data point is represented as a node in a graph. In comparison, TabularNet (Du et al. 2021) and Table2Graph (Zhou et al. 2022) follow the second approach, where the first uses a Graph Convolutional Network to model the relationships between features, and the second learns a probability adjacency matrix for a unified graph that models the interaction between features of the data points. GNNs can also be combined with other algorithms suited for tabular data, e.g., as in BGNN (Ivanov and Prokhorenkova 2021), which combines gradient-boosted decision trees and a GNN in one pipeline, where the GNN addresses the graph structure and the gradient-boosted decision trees handle the heterogeneous features of the tabular data. To the best of our knowledge, all previous approaches to using GNNs for tabular data result in black-box models and they are hence associated with the issues discussed above when applied in contexts with strong requirements on trustworthiness. In this work, we propose a novel GNN approach for tabular data, with the aim to eliminate the need to apply post-hoc explanation techniques without sacrificing predictive performance.

The main contributions of this study are:

- a novel approach, called **I**nterpretable **G**raph **N**eural **Net**work for tabular data (IGNNet), that exploits powerful graph neural network models while still being able to show exactly how the prediction is derived from the input features in a transparent way

- a large-scale empirical investigation evaluating the explanations of IGNNet as well as comparing the predictive performance of IGNNet to state-of-the-art approaches for tabular data; XGBoost, Random Forests, and multi-layer

perceptron (MLP), as well as to an algorithm generating interpretable models; TabNet (Arik and Pfister 2021)

- an ablation study comparing the performance of the proposed approach to a black-box version, i.e., not constraining the learning algorithm to produce transparent models for the predictions

In the next section, we briefly review related work. In Section 3, we describe the proposed interpretable graph neural network. In Section 4, results from a large-scale empirical investigation are presented and discussed, in which the explanations of the proposed method are evaluated and the performance is compared both to interpretable and powerful black-box models. Finally, in the concluding remarks section, we summarize the main conclusions and point out directions for future work.

## 2 Related Work

In this section, we provide some pointers to self-explaining (regular and graph) neural networks, and briefly discuss their relation to model-agnostic explanation techniques and interpretable models. We also provide pointers to work on interpretable deep learning approaches for tabular data.

### 2.1 Self-Explaining Neural Networks

Several approaches to generating so-called self-explaining neural networks have been introduced in the literature; in addition to generating a prediction model, they all incorporate a component for explaining the predictions. They can be seen as model-specific explanation techniques, in contrast to model-agnostic techniques, such as LIME and SHAP, but sharing the same issues regarding fidelity and lack of detail regarding the exact computation of the predictions.

Approaches in this category include the method in (Lei, Barzilay, and Jaakkola 2016), which is an early self-explaining neural network for text classification, the Contextual Explanation Network (CEN) (Al-Shedivat, Dubey, and Xing 2022), which generates explanations using intermediate graphical models, the Self-explaining Neural Network (SENN) (Alvarez Melis and Jaakkola 2018), which generalizes linear classifiers to neural networks using a concept autoencoder, and the CBM-AUC (Sawada and Nakamura 2022), which improves the efficiency of the former by replacing the decoder with a discriminator. Some approaches generate explanations in the form of counterfactual examples, e.g., CounterNet (Guo, Nguyen, and Yadav 2021), and Variational Counter Net (VCNet) (Guyomard, Fessant, and Guyet 2022). Again, such explanations do not provide detailed information on how the original predictions are computed and how exactly the input features affect the outcome.

### 2.2 Self-Explaining Graph Neural Networks

The Self-Explaining GNN (SE-GNN) (Dai and Wang 2021) uses similarities between nodes to make predictions on the nodes' labels and provide explanations using the most similar K nodes with labels. ProtGNN (Zhang et al. 2022) also computes similarities, but between the input graph and prototypical graph patterns that are learned per class. Cui et

al. (Cui et al. 2022) proposed a framework to build interpretable GNNs for connectome-based brain disorder analysis that resembles the signal correlation between different brain areas.

Similar to the (standard) self-explaining neural networks, the approaches that target graph neural networks provide abstract views of the predictions; the users cannot trace the exact computations, in contrast to when using interpretable models, which in principle allows for the inferences to be executed by hand.

### 2.3 Interpretable Deep Learning for Tabular Data

In an endeavor to provide an interpretable regression model for tabular data while retaining the performance of deep learning models, and inspired by generalized linear models (GLM), LocalGLMnet was proposed to make the regression parameters of a GLM feature dependent, allowing for quantifying variable importance and also conducting variable selection (Richman and Wüthrich 2022). TabNet (Arik and Pfister 2021) is another interpretable method proposed for tabular data learning, which employs a sequential attention mechanism and learnable masks for selecting a subset of meaningful features to reason from at each decision step. The feature selection is instance-based, i.e., it differs from one instance to another. The feature selection masks can be visualized to highlight important features and show how they are combined. However, it is not obvious how the features are actually used to form the predictions.

## 3 The Proposed Approach: IGNNet

This section describes the proposed method to produce an interpretable model using a graph neural network. We first outline the details of a GNN for graph classification and then show how it can be constrained to produce interpretable models. Afterward, we show how it can be applied to tabular data. Finally, we show how the proposed approach can maintain both interpretability and high performance.

### 3.1 Interpretable Graph Neural Network

The input to a GNN learning algorithm is a set of graphs denoted by $\mathcal{G} = (V, E, X, \mathcal{W})$, consisting of a set of nodes $V$, a set of edges $E$, a set of node feature vectors $X$, and a set of edge weights $\mathcal{W}$, where $V = \{v_1, \ldots, v_N\}$, $E \subseteq \{(v_i, v_j) | v_i, v_j \in V\}$, $X = \{\mathbf{x}_1, \ldots, \mathbf{x}_N\}$, and the weight of edge $(v_i, v_j)$ is represented by a scalar value $\delta_{i,j}$ in the set of edge weights $\mathcal{W}$, where $\delta_{i,j} = \mathcal{W}(i, j)$. A GNN algorithm learns a representation vector $\mathbf{h}_i$ for each node $v_i$, which is initialized as $\mathbf{h}_i^{(0)} = \mathbf{x}_i$. The key steps in a GNN for graph classification can be summarized by the following two phases (Xu et al. 2019):

(a) **Message Passing:** Each node passes a message to the neighboring nodes, then aggregates the passed information from the neighbors. Finally, the node representation is updated with the aggregated information. A neural network can also be used to learn some message functions between nodes. The message passing phase can be formulated as:

$$\mathbf{h}_i^{(l+1)} = \varphi(\mathbf{w}^{(l+1)}(\sum_{u \in \mathcal{N}(i)} \delta_{i,u}\mathbf{h}_u^{(l)} + \delta_{i,i}\mathbf{h}_i^{(l)})) \quad (1)$$

where $\delta_{i,u}$ is the weight assigned to the edge between node $v_i$ and node $v_u$. $\mathbf{h}_i^{(l)}$ is the hidden representation of the node $v_i$ in the $l$-th layer, $\mathbf{w}^{(l+1)}$ represents the learnable parameters, and $\varphi$ is a non-linearity function.

The adjacency matrix $\boldsymbol{A}$ of size $|V| \times |V|$ contains the edge weights and can be normalized similar to a Graph Convolutional Network (GCN) (Kipf and Welling 2017) as shown in (2).

$$\tilde{\boldsymbol{A}} = \boldsymbol{D}^{-\frac{1}{2}}\boldsymbol{A}\boldsymbol{D}^{-\frac{1}{2}} \quad (2)$$

Here $\boldsymbol{D}$ is the degree matrix $\boldsymbol{D}_{ii} = \sum_j \boldsymbol{A}_{ij}$ (Kipf and Welling 2017).

(b) **Graph Pooling (Readout):** A representation of the whole graph $\mathcal{G}$ is learned using a simple or advanced function (Xu et al. 2019), e.g, sum, mean, or MLP.

The whole graph representation obtained from the **graph pooling** phase can be submitted to a classifier to predict the class of the graph, which can be trained in an end-to-end architecture (Zhang et al. 2018; Ying et al. 2018).

*The pooling function can be designed to provide an interpretable graph classification layer.* Thus, the final hidden representation of each node is mapped to a single value, for instance, through a neural network layer or dot product ($f(\mathbf{h}_i^{(l+1)} \in \mathbb{R}^n) = h_i \in \mathbb{R}^1$), and concatenated to obtain the final representation $\mathbf{g}$ of the graph $\mathcal{G}$ where a scalar value in $\mathbf{g}$ corresponds to a node in the graph. Consequently, if a set of weights is applied to classify the graph, we can trace the contribution of each node to the predicted outcome, i.e., the user can find out which nodes contributed to the predicted class. For example, $\mathbf{g}$ can be used directly as follows:

$$\hat{y} = \text{link}(\sum_{i=1}^{n} \mathbf{w}_i\mathbf{g}_i) \quad (3)$$

where $\mathbf{w}_i$ is the weight assigned to node $v_i$ represented in $\mathbf{g}_i$. The link function is applied to accommodate a valid range of outputs, e.g., the sigmoid function for binary and softmax for multi-class classification. This is equivalent to:

$$\hat{y} = \text{link}(\sum_{i=1}^{n} \mathbf{w}_i f(\mathbf{h}_i^{(l+1)})) \quad (4)$$

In the case of binary classification, one vector of weights ($\mathbf{w}$) is applied, and for multiple classes, each class has a separate vector of weights.

### 3.2 Representing Tabular Data Points as Graphs

The proposed readout function in the previous subsection allows for determining the contribution of each node in a prediction, if a white-box classification layer is used for the latter. Therefore, we propose representing each data instance as a graph where *the features are the nodes* of that graph and *the linear correlation between features are the edge weights*,

as we assume that not all features are completely independent. The initial representation of a node is a vector of one dimension, and the value is just the feature value, which can be embedded into a higher dimensionality. The idea is outlined in Algorithm 1 and illustrated in Figure 1.

---

**Algorithm 1:** IGNNet

**Data:** a set of graphs $\mathbb{G}$ and labels $y$
**Result:** Model parameters $\theta$
Initialize $\theta$
**for** *number of training iterations* **do**
    $\mathcal{L} \leftarrow 0$
    **for** *each $\mathcal{G}_j \in \mathbb{G}$* **do**
        $H_j \leftarrow$ messagePassing($\mathcal{G}_j$)
        $\mathbf{g}_j \leftarrow$ readout($H_j$)
        $\hat{y}_j \leftarrow$ predict($\mathbf{g}_j$)
        $\mathcal{L} \leftarrow \mathcal{L} + \text{loss}(\hat{y}_j, y_j)$
    **end**
    Compute gradients $\nabla_\theta \mathcal{L}$
    Update $\theta \leftarrow \theta - \nabla_\theta \mathcal{L}$
**end**

---

In order to make a prediction for a test instance, the data point has to be converted to a graph using the same procedure for building the input graphs to IGNNet, and for which graph node representations are obtained using a GNN with parameters $\theta$. Finally, the output layer is used to form the prediction.

### 3.3 How can IGNNet achieve high performance while maintaining interpretability?

An expressive GNN can potentially capture complex patterns and dependencies in the graph, allowing nodes to be mapped to distinct representations based on their characteristics and relationships (Li and Leskovec 2022). Moreover, a GNN with an injective aggregation scheme can not only distinguish different structures but also map similar structures to similar representations (Xu et al. 2019). Therefore, if the tabular data are properly presented as graphs, GNNs with the aforementioned expressive capacities can model relationships and interactions between features, and consequently approximate complex non-linear mappings from inputs to predictions. On top of that, it has been shown by (ENNADIR et al. 2023) that GCNs based on 1-Lipschitz continuous activation functions can be improved in stability and robustness with Lipschitz normalization and continuity analysis; similar findings have also been demonstrated on graph attention networks (GAT) (Dasoulas, Scaman, and Virmaux 2021). This property is of particular importance when the application domain endures adversarial attacks or incomplete tabular data.

The proposed readout function in subsection 3.1 can produce an interpretable output layer. However, it does not guarantee the interpretability of the whole GNN without message-passing layers that consistently maintain relevant representations of the input features. Accordingly, we constrain the message-passing layer to produce interpretable
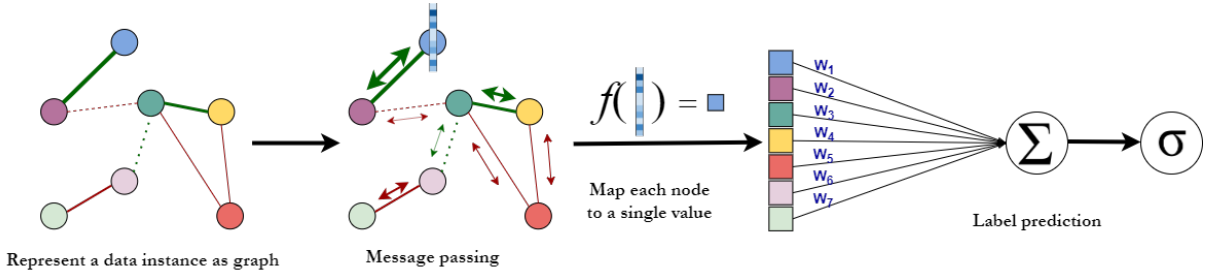
Figure 1: **An overview of our proposed approach.** Each data instance is represented as a graph by embedding the feature values into a higher dimensionality, and the edge between two features (nodes) is the correlation value. Multiple iterations of message passing are then applied. Finally, the learned node representation is projected into a single value, and a whole graph representation is obtained by concatenating the projected values.

models using the following conditions:

1. Each feature is represented in a distinct node throughout the consecutive layers.
2. Each node is bounded to interact with a particular neighborhood, where it maintains correlations with the nodes within that neighborhood. As a result, the aggregated messages could potentially hold significance to the input feature values.

Since the edge weights are the correlation values, they determine the strength and the sign of the messages obtained from a neighborhood, allowing each node to store information not only concerning the original input feature but also features that are correlated with it. The proposed graph pooling function, combined with the constrained message-passing layers that keep representative information about the input features, allows tracking each feature's contribution at the output layer and also through the message-passing layers all the way to the input features.

## 4 Empirical Investigation

This section evaluates both the explanations and predictive performance of IGNNet. We begin by outlining the experimental setup, then by evaluating the explanations produced by IGNNet, and lastly, we benchmark the predictive performance.

### 4.1 Experimental Setup

A GNN consists of one or more message-passing layers, and each layer can have a different design, e.g., different activation functions and batch normalization, which is the intra-layer design level (You, Ying, and Leskovec 2020). There is also the inter-layer design level which involves, for instance, how the message-passing layers are organized into a neural network and if the skip connections are added between layers (You, Ying, and Leskovec 2020). While the intra-layer and inter-layer designs can vary based on the nature of the prediction task, we propose a general architecture for our empirical investigation. However, it is up to the user to modify the architecture of the GNN. The number of message-passing layers, number of units in linear transformations, and other hyperparameters were found based on a quasi-random search and evaluation on development sets of the following three datasets: Churn, Electricity, and Higgs. We have six message-passing layers in the proposed architecture, each with a Relu activation function. Multiple learnable weights are also applied to the nodes' representation, followed by a Relu function. Besides three batch normalization layers, four skip connections are added as illustrated in Figure 2. After all the GNN layers, we use a feedforward neural network (FNN) to map the multidimensional representation of each node into a single value. In the FNN, we do not include any activation functions in order to keep the mapping linear, but a sigmoid function is applied after the final layer to obtain a value between 0 and 1 for each node. The FNN is composed of 8 layers with the following numbers of units (128, 64, 32, 16, 8, 4, 2, 1) and 3 batch normalization layers after the second, fourth, and sixth hidden layers. After the FNN, the nodes' final values are concatenated to form a representation of the whole graph (data instance). Finally, the weights that are output are used to make predictions. The GNN is trained end-to-end, starting from the embeddings layer and ending with the class prediction.

We also provide an opaque variant of IGNNet (OGNNet, opaque graph neural net for tabular data), where the FNN, along with the output layer, is replaced by an MLP of one hidden layer with 1024 hidden units and a Relu activation function. All the learned node representations just before the FNN are concatenated and passed to the MLP for class prediction. The OGNNet is introduced to determine, by an ablation study, how much predictive performance we may lose by squashing the learned multidimensional representation of the nodes into scalar values and applying a white box classifier instead of a black box.

In the experiments, 35 publicly available datasets are used.[1] Each dataset is split into training, development, and test sets. The development set is used for overfitting detection and early stopping of the training process, the training set is used to train the model, and the test set is used to evaluate the model.[2] For a fair comparison, all the compared learning algorithms are trained without hyperparam-

---

[1] All the datasets were obtained from https://www.openml.org

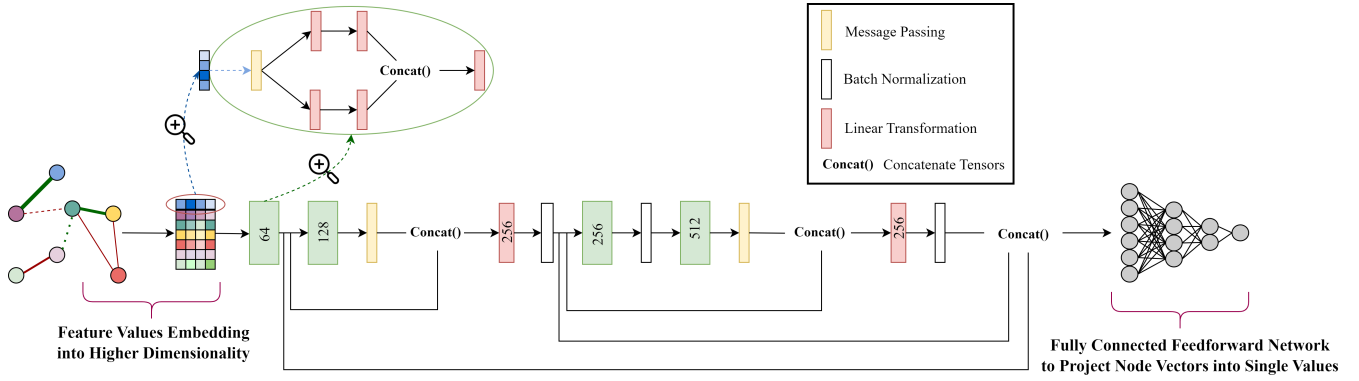[2] Detailed information about each dataset is provided in the appendix.

Figure 2: **IGNNet default architecture.** It starts with the embedding layer, a linear transformation from one dimension to 64 dimensions. A Relu activation function follows each message-passing layer and each green block as well. The feedforward network at the end has no activation functions between layers to ensure a linear transformation into a single value. A sigmoid activation function follows the feedforward network to obtain the final value for each feature between 0 and 1.

eters tuning using the default settings on each dataset. In cases where the learning algorithm does not employ the development set to measure performance progress for early stopping, the development and training subsets are combined into a joint training set. The adjacency matrix uses the correlation values computed on the training data split. The weight on edge from the node to itself (self-loop) is a user-adjustable hyperparameter, constitutes between 70% to 90% (on average) of the weighted summation to keep a strong message per node that does not fade out with multiple layers of message-passing. Weak correlation values are excluded from the graph, so if the absolute correlation value is below 0.2, the edge is removed unless no correlation values are above 0.2; in case of the latter, the procedure is repeated using a reduced threshold of 0.05.[3] The Pearson correlation coefficient (Pearson 1895) is used to estimate the linear relationship between features. In the data preprocessing step, the categorical features are binarized using one-hot encoding, and all the feature values are normalized using min-max normalization (the max and min values are computed on the training split). The normalization keeps the feature values between 0 and 1, which is essential for the IGNNet to have one scale for all nodes.

The following algorithms are also evaluated in the experiments: XGBoost, Random Forests, MLP and TabNet. XGBoost and Random Forests are trained on the combined training and development sets. The MLP has two layers of 1024 units with Relu activation function and is trained on the combined training and development sets with early stopping and 0.1 validation fraction. TabNet is trained with early stopping after 20 consecutive epochs without improvement on the development set, and the best model is used in the evaluation.

For imbalanced binary classification datasets, we oversample the minority class in the training set to align the size with the majority class. All the compared algorithms

are trained using the oversampled training data. While for multi-class datasets, no oversampling is conducted. The area under the ROC curve (AUC) is used to measure the predictive performance, as it is not affected by the decision threshold. For the multi-class datasets, weighted AUC is calculated, i.e., the AUC is computed for each class against the rest and weighted by the support.

### 4.2 Evaluation of Explanations

The feature scores produced by IGNNet should ideally reflect the contribution of each feature toward the predicted outcome and, therefore, they should be equivalent to the true Shapley values. As it has been shown that KernelSHAP converges to the true Shapley values when provided with an infinite number of samples (Covert and Lee 2021; Jethani et al. 2022), it is anticipated that the explanations generated by KernelSHAP will progressively converge to more similar values to the scores of IGNNet as the sampling process continues. This convergence arises from KernelSHAP moving towards the true values, while the scores of IGNNet are expected to align with these true values. To examine this conjecture, we explain IGNNet using KernelSHAP and measure the similarity between KernelSHAP's explanations and IGNNet's scores following each iteration of data sampling and KernelSHAP evaluation.[4] For the feasibility of the experiment, 500 examples are randomly selected from the test set of each dataset to be explained. The cosine similarity and Spearman rank-order correlation are used to quantify the similarity between explanations. The cosine similarity measures the similarity in the orientation (Han, Kamber, and Pei 2012), while the Spearman rank-order measures the similarity in ranking the importance scores (Rahnama et al. 2021).

The results demonstrate a general trend wherein KernelSHAP's explanations converge to more similar values to IGNNet's scores across various data instances and the 35

---

[3]In the technical appendix, we provide an ablation study on the effects of different hyperparameter preferences on predictive performance.

[4]KernelSHAP experiments were conducted using the following open-source implementation: https://github.com/iancovert/shapley-regression
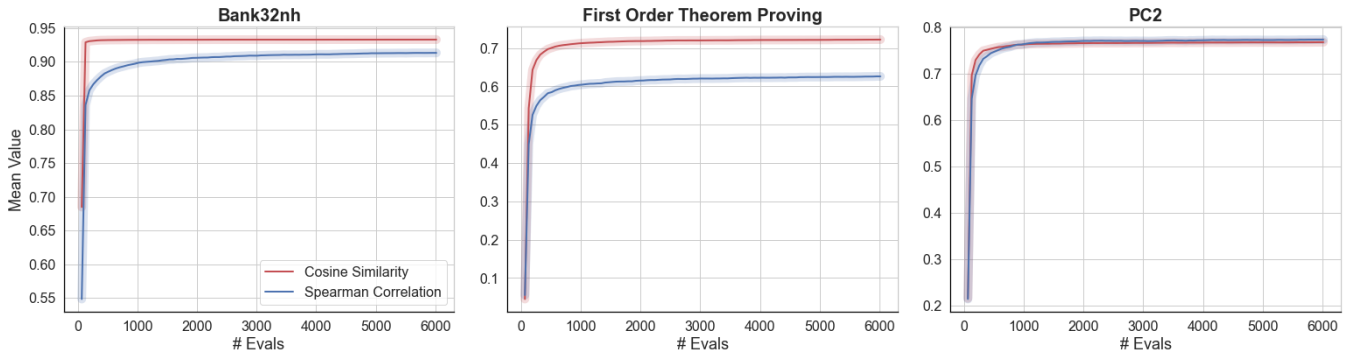
Figure 3: **Comparison of KernelSHAP's approximations and the importance scores obtained from IGNNet.** We measure the similarity of KernelSHAP's approximations to the scores of IGNNet at each iteration of data sampling and evaluation of KernelSHAP. KernelSHAP exhibits improvement in approximating the scores derived from IGNNet with more data sampling.

datasets, as depicted in Figure 10. The consistent convergence to more similar values clearly indicates that IGNNet provides transparent models with feature scores aligned with the true Shapley values.[5][6]

## 4.3 Evaluation of Predictive Performance

Detailed results for IGNNet and the five competing algorithms on the 35 datasets are shown in Table 1. The ranking of the five algorithms across 35 datasets, based on their AUC values, reveals OGNNet to exhibit superior performance, claiming the top position, closely followed by IGN-Net and XGBoost. In order to investigate whether the observed differences are statistically significant, the Friedman test (Friedman 1939) was employed, which indeed allowed to reject the null hypothesis, i.e., that there is no difference in predictive performance, as measured by the AUC, at the 0.05 level. The result of subsequently applying the post-hoc Nemenyi test (Nemenyi 1963) to determine what pairwise differences are significant, again at the 0.05 level, is summarized in Figure 4. However, the result shows no specific significant pairwise differences between any of the compared algorithms. Furthermore, the results show that using IGNNet instead of the black-box variant, OGNNet, does not significantly reduce the predictive performance while maintaining performance at the level of other powerful algorithms for tabular data, e.g., XGBoost and Random Forests.

## 5 Concluding Remarks

We have proposed IGNNet, an algorithm for tabular data classification, which exploits graph neural networks to produce transparent models. In contrast to post-hoc explanation techniques, IGNNet does not approximate or require costly computations, but provides the explanation while computing the prediction, and where the explanation prescribes exactly how the prediction is computed.

---

[5]The complete results of the 35 datasets are provided in the technical appendix.

[6]We provide illustrations of individual explanations in the technical appendix.

We have presented results from a large-scale empirical investigation, in which IGNNet was evaluated with respect to explainability and predictive performance. IGNNet was shown to generate explanations with feature scores aligned with the Shapley values without further computational cost. IGNNet was also shown to achieve a similar predictive performance as XGBoost, Random Forests, TabNet, and MLP, which are all well-known for their ability to generate high-performing models.

One direction for future research is to explore approaches to model feature interaction in the adjacency matrix that go beyond linear correlations. Understanding how such non-linear interactions between features may impact the model's interpretability could be an intriguing area of exploration. A second direction is to investigate alternative encoders for categorical features rather than relying on one-hot encoding. It would also be interesting to extend IGNNet to handle non-tabular datasets, including images and text, which would require entirely different approaches to representing each data point as a graph. Another important direction for future work is to use IGNNet for studying possible adversarial attacks on a predictive model. Finally, an important direction would be to complement the empirical evaluation with user-grounded evaluations, e.g., measuring to what extent certain tasks could be more effectively and efficiently solved when the users are provided with a transparent model that shows how the prediction has been computed from the input.

## 6 Acknowledgement

## References

Al-Shedivat, M.; Dubey, A.; and Xing, E. 2022. Contextual Explanation Networks. *J. Mach. Learn. Res.*, 21(1).

Alvarez Melis, D.; and Jaakkola, T. 2018. Towards Robust Interpretability with Self-Explaining Neural Networks. In Bengio, S.; Wallach, H.; Larochelle, H.; Grauman, K.; Cesa-

```
CD
|————————|
1        2        3        4        5        6

OGNNet ————————————————————                                    TabNet
IGNNet ————————————————————                                    MLP
XGBoost ————————————————————                                   Random Forests
```
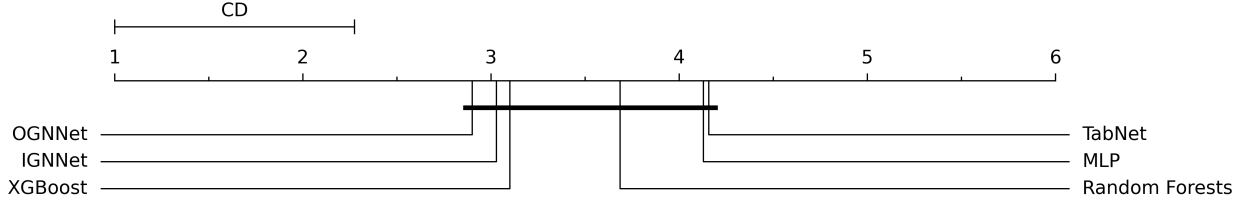
Figure 4: **The average rank of the compared classifiers on the 35 datasets with respect to the AUC** (a lower rank is better), where the critical difference (CD) represents the largest difference that is not statistically significant.

Table 1: The AUC of IGNNet, OGNNet, MLP, Random Forests, and XGBoost. The best-performing model is colored in blue , and the second best-performing is colored in light blue .

| Dataset | IGNNet | OGNNet | TabNet | MLP | Random Forests | XGBoost |
|---|---|---|---|---|---|---|
| Abalone | 0.881 | 0.883 | 0.857 | 0.877 | 0.876 | 0.869 |
| Ada Prior | 0.905 | 0.888 | 0.848 | 0.877 | 0.885 | 0.894 |
| Adult | 0.917 | 0.915 | 0.919 | 0.881 | 0.907 | 0.931 |
| Bank 32 nh | 0.887 | 0.887 | 0.881 | 0.859 | 0.876 | 0.874 |
| Covertype | 0.984 | 0.988 | 0.969 | 0.861 | 0.995 | 0.967 |
| Credit Card Fraud | 0.987 | 0.966 | 0.969 | 0.913 | 0.914 | 0.975 |
| Delta Ailerons | 0.977 | 0.977 | 0.974 | 0.977 | 0.978 | 0.977 |
| Electricity | 0.901 | 0.929 | 0.894 | 0.928 | 0.97 | 0.973 |
| Elevators | 0.951 | 0.948 | 0.95 | 0.95 | 0.913 | 0.943 |
| HPC Job Scheduling | 0.908 | 0.921 | 0.775 | 0.908 | 0.952 | 0.955 |
| Fars | 0.956 | 0.958 | 0.954 | 0.955 | 0.949 | 0.962 |
| 1st Order Theorem Proving | 0.776 | 0.808 | 0.495 | 0.805 | 0.854 | 0.858 |
| Helena | 0.875 | 0.889 | 0.884 | 0.897 | 0.855 | 0.875 |
| Heloc | 0.783 | 0.787 | 0.772 | 0.783 | 0.778 | 0.775 |
| Higgs | 0.762 | 0.785 | 0.804 | 0.774 | 0.793 | 0.797 |
| Indian Pines | 0.984 | 0.992 | 0.99 | 0.973 | 0.979 | 0.987 |
| Jannis | 0.856 | 0.857 | 0.867 | 0.861 | 0.861 | 0.872 |
| JM1 | 0.739 | 0.725 | 0.711 | 0.728 | 0.747 | 0.733 |
| LHC Identify Jets | 0.941 | 0.941 | 0.944 | 0.861 | 0.935 | 0.941 |
| Madelon | 0.906 | 0.718 | 0.501 | 0.668 | 0.79 | 0.891 |
| Magic Telescope | 0.907 | 0.921 | 0.927 | 0.929 | 0.934 | 0.928 |
| MC1 | 0.957 | 0.904 | 0.89 | 0.853 | 0.844 | 0.943 |
| Mozilla4 | 0.954 | 0.961 | 0.971 | 0.963 | 0.988 | 0.99 |
| Microaggregation2 | 0.778 | 0.792 | 0.752 | 0.782 | 0.768 | 0.781 |
| Numerai28.6 | 0.526 | 0.534 | 0.52 | 0.518 | 0.519 | 0.514 |
| Otto Group Product | 0.96 | 0.971 | 0.968 | 0.972 | 0.973 | 0.974 |
| PC2 | 0.881 | 0.815 | 0.844 | 0.571 | 0.55 | 0.739 |
| Phonemes | 0.922 | 0.96 | 0.898 | 0.93 | 0.964 | 0.953 |
| Pollen | 0.492 | 0.508 | 0.509 | 0.496 | 0.489 | 0.467 |
| Satellite | 0.998 | 0.993 | 0.911 | 0.992 | 0.998 | 0.992 |
| Scene | 0.994 | 0.989 | 0.986 | 0.992 | 0.983 | 0.982 |
| Speed Dating | 0.853 | 0.835 | 0.797 | 0.822 | 0.845 | 0.86 |
| Telco Customer Churn | 0.858 | 0.845 | 0.841 | 0.783 | 0.84 | 0.843 |
| Vehicle sensIT | 0.918 | 0.918 | 0.917 | 0.914 | 0.912 | 0.916 |
| Waveform-5000 | 0.965 | 0.962 | 0.933 | 0.965 | 0.959 | 0.957 |

Bianchi, N.; and Garnett, R., eds., *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc.

Arik, S. O.; and Pfister, T. 2021. TabNet: Attentive Interpretable Tabular Learning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(8): 6679–6687.

Berkson, J. 1944. Application of the Logistic Function to Bio-Assay. *Journal of the American Statistical Association*, 39(227): 357–365.

Breiman, L. 2001. Random Forests. *Machine Learning*, 45(1): 5–32.

Chen, T.; and Guestrin, C. 2016. XGBoost: A Scalable Tree Boosting System. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '16, 785–794. New York, NY, USA: Association for Computing Machinery. ISBN 9781450342322.

Covert, I.; and Lee, S.-I. 2021. Improving KernelSHAP: Practical Shapley Value Estimation Using Linear Regression. In Banerjee, A.; and Fukumizu, K., eds., *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*, volume 130 of *Proceedings of Machine Learning Research*, 3457–3465. PMLR.

Cui, H.; Dai, W.; Zhu, Y.; Li, X.; He, L.; and Yang, C. 2022. Interpretable Graph Neural Networks for Connectome-Based Brain Disorder Analysis. In *Medical Image Computing and Computer Assisted Intervention – MICCAI 2022: 25th International Conference, Singapore, September 18–22, 2022, Proceedings, Part VIII*, 375–385. Berlin, Heidelberg: Springer-Verlag. ISBN 978-3-031-16451-4.

Dai, E.; and Wang, S. 2021. Towards Self-Explainable Graph Neural Network. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, CIKM '21, 302–311. New York, NY, USA: Association for Computing Machinery. ISBN 9781450384469.

Dasoulas, G.; Scaman, K.; and Virmaux, A. 2021. Lipschitz normalization for self-attention layers with application to graph neural networks. In Meila, M.; and Zhang, T., eds., *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, 2456–2466. PMLR.

Delaunay, J.; Galárraga, L.; and Largouët, C. 2020. Improving Anchor-based Explanations. In *CIKM 2020 - 29th ACM International Conference on Information and Knowledge Management*, 3269–3272. Galway / Virtual, Ireland: ACM.

Du, L.; Gao, F.; Chen, X.; Jia, R.; Wang, J.; Zhang, J.; Han, S.; and Zhang, D. 2021. TabularNet: A Neural Network Architecture for Understanding Semantic Structures of Tabular Data. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining (KDD'21)*.

ENNADIR, S.; ABBAHADDOU, Y.; Vazirgiannis, M.; and Boström, H. 2023. A Simple and Yet Fairly Effective Defense for Graph Neural Networks. In *The Second Workshop on New Frontiers in Adversarial Machine Learning*.

Friedman, M. 1939. A Correction: The use of ranks to avoid the assumption of normality implicit in the analysis of variance. *Journal of the American Statistical Association*, 34(205): 109–109.

Goodman, B.; and Flaxman, S. 2017. European Union Regulations on Algorithmic Decision-Making and a "Right to Explanation". *AI Magazine*, 38(3): 50–57.

Guo, H.; Nguyen, T.; and Yadav, A. 2021. CounterNet: End-to-End Training of Counterfactual Aware Predictions. In *ICML 2021 Workshop on Algorithmic Recourse*.

Guo, X.; Quan, Y.; Zhao, H.; Yao, Q.; Li, Y.; and Tu, W. 2021. TabGNN: Multiplex Graph Neural Network for Tabular Data Prediction. *CoRR*, abs/2108.09127.

Guyomard, V.; Fessant, F.; and Guyet, T. 2022. VCNet: A self-explaining model for realistic counterfactual generation. In *Machine Learning and Knowledge Discovery in Databases - European Conference, ECML PKDD 2022, Grenoble, France, September 19-23*.

Han, J.; Kamber, M.; and Pei, J. 2012. 2 - Getting to Know Your Data. In Han, J.; Kamber, M.; and Pei, J., eds., *Data Mining (Third Edition)*, The Morgan Kaufmann Series in Data Management Systems, 39–82. Boston: Morgan Kaufmann, third edition edition. ISBN 978-0-12-381479-1.

Ivanov, S.; and Prokhorenkova, L. 2021. Boost then Convolve: Gradient Boosting Meets Graph Neural Networks. In *International Conference on Learning Representations (ICLR)*.

Jethani, N.; Sudarshan, M.; Covert, I. C.; Lee, S.-I.; and Ranganath, R. 2022. FastSHAP: Real-Time Shapley Value Estimation. In *International Conference on Learning Representations*.

Kipf, T. N.; and Welling, M. 2017. Semi-Supervised Classification with Graph Convolutional Networks. In *International Conference on Learning Representations*.

Lakkaraju, H.; Kamar, E.; Caruana, R.; and Leskovec, J. 2017. Interpretable & Explorable Approximations of Black Box Models. *CoRR*, abs/1707.01154.

Lei, T.; Barzilay, R.; and Jaakkola, T. 2016. Rationalizing Neural Predictions. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, 107–117. Austin, Texas: Association for Computational Linguistics.

Li, P.; and Leskovec, J. 2022. The Expressive Power of Graph Neural Networks. In Wu, L.; Cui, P.; Pei, J.; and Zhao, L., eds., *Graph Neural Networks: Foundations, Frontiers, and Applications*, 63–98. Singapore: Springer Singapore.

Lundberg, S. M.; and Lee, S.-I. 2017. A Unified Approach to Interpreting Model Predictions. In Guyon, I.; Luxburg, U. V.; Bengio, S.; Wallach, H.; Fergus, R.; Vishwanathan, S.; and Garnett, R., eds., *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.

Mori, T.; and Uchihira, N. 2019. Balancing the trade-off between accuracy and interpretability in software defect prediction. *Empirical Software Engineering*, 24(2): 779–825.

Nemenyi, P. B. 1963. *Distribution-free multiple comparisons*. Ph.D. thesis, Princeton University.

Pearson, K. 1895. Note on Regression and Inheritance in the Case of Two Parents. *Proceedings of the Royal Society of London Series I*, 58: 240–242.

Pintelas, E.; Livieris, I. E.; and Pintelas, P. 2020. A Grey-Box Ensemble Model Exploiting Black-Box Accuracy and White-Box Intrinsic Interpretability. *Algorithms*, 13(1).

Rahnama, A. H. A.; Bütepage, J.; Geurts, P.; and Boström, H. 2021. Evaluation of Local Model-Agnostic Explanations Using Ground Truth. *CoRR*, abs/2106.02488.

Ribeiro, M. T.; Singh, S.; and Guestrin, C. 2016. "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, San Francisco, CA, USA, August 13-17, 2016*, 1135–1144.

Ribeiro, M. T.; Singh, S.; and Guestrin, C. 2018. Anchors: High-Precision Model-Agnostic Explanations. In *AAAI Conference on Artificial Intelligence (AAAI)*.

Richman, R.; and Wüthrich, M. V. 2022. LocalGLMnet: interpretable deep learning for tabular data. *Scandinavian Actuarial Journal*, 0(0): 1–25.

Rudin, C. 2019. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature Machine Intelligence*, 1(5): 206–215.

Sawada, Y.; and Nakamura, K. 2022. Concept Bottleneck Model With Additional Unsupervised Concepts. *IEEE Access*, 10: 41758–41765.

Wilcoxon, F. 1945. Individual comparisons by ranking methods. biometrics bulletin 1, 6 (1945), 80–83.

Xu, K.; Hu, W.; Leskovec, J.; and Jegelka, S. 2019. How Powerful are Graph Neural Networks? In *International Conference on Learning Representations*.

Yeh, C.-K.; Hsieh, C.-Y.; Suggala, A.; Inouye, D. I.; and Ravikumar, P. K. 2019. On the (In)fidelity and Sensitivity of Explanations. In Wallach, H.; Larochelle, H.; Beygelzimer, A.; d'Alché-Buc, F.; Fox, E.; and Garnett, R., eds., *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.

Ying, R.; You, J.; Morris, C.; Ren, X.; Hamilton, W. L.; and Leskovec, J. 2018. Hierarchical Graph Representation Learning with Differentiable Pooling. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, NIPS'18, 4805–4815. Red Hook, NY, USA: Curran Associates Inc.

You, J.; Ying, R.; and Leskovec, J. 2020. Design Space for Graph Neural Networks. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS'20. Red Hook, NY, USA: Curran Associates Inc. ISBN 9781713829546.

Zhang, M.; Cui, Z.; Neumann, M.; and Chen, Y. 2018. An End-to-End Deep Learning Architecture for Graph Classification. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence and Thirtieth Innovative Applications of Artificial Intelligence Conference and Eighth AAAI Symposium on Educational Advances in Artificial Intelligence*, AAAI'18/IAAI'18/EAAI'18. AAAI Press. ISBN 978-1-57735-800-8.

Zhang, Z.; Liu, Q.; Wang, H.; Lu, C.; and Lee, C.-K. 2022. ProtGNN: Towards Self-Explaining Graph Neural Networks. In *AAAI*.

Zhou, K.; Liu, Z.; Chen, R.; Li, L.; Choi, S.-H.; and Hu, X. 2022. Table2Graph: Transforming Tabular Data to Unified Weighted Graph. In Raedt, L. D., ed., *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22*, 2420–2426. International Joint Conferences on Artificial Intelligence Organization. Main Track.
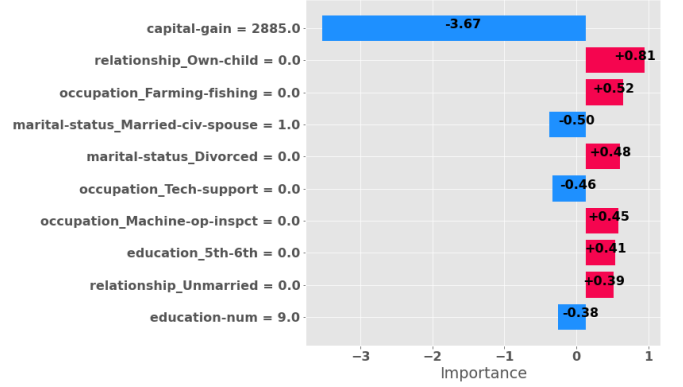
# Appendix

## Illustration of Explanations

In this section, we show how the computed feature scores by IGNNet can be used to understand the feature contributions toward a specific prediction. Note that this is done in exactly the same way as how one would interpret the predictions of a logistic regression model or the feature importance scores generated by the SHAP explainer (Lundberg and Lee 2017). As the computed feature scores reveal exactly how IGNNet formed the prediction, the user can directly see which features have the greatest impact on the final prediction, and possibly also how they may be modified to affect the outcome. To demonstrate this, we present the feature scores for predictions made by IGNNet using two examples from the Adult dataset and the Churn dataset. In the following illustrations, we display the feature scores centered around the bias value, which, when summed with the bias, will produce the exact outcome of IGNNet if the sigmoid function is applied. The scores are sorted according to their absolute values, and only the top 10 features are plotted for ease of presentation. A displayed score $\tau_i$ of feature $x_i$ represents all the weights and the computations applied to the input value, as shown in equation 5.
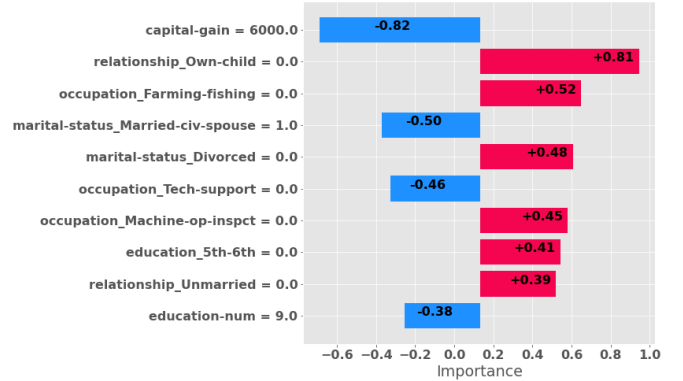
$$
\begin{aligned}
\tau_i = & \mathbf{w}_i f\big(\varphi(\mathbf{w}^{(l+1)}\big(\sum_{u \in \mathcal{N}(i)} \delta_{i,u} \mathbf{h}_u^{(l)} \\
& + \delta_{i,i} \varphi(\mathbf{w}^{(l)}\big(\sum_{u \in \mathcal{N}(i)} \delta_{i,u} \mathbf{h}_u^{(l-1)} \\
& + \delta_{i,i}(\dots \varphi(\mathbf{w}^{(1)}\big(\sum_{u \in \mathcal{N}(i)} \delta_{i,u}\mathbf{x}_u + \delta_{i,i}\mathbf{x}_i)\dots))))))
\end{aligned}
\tag{5}
$$

The first example, from the Adult dataset, shown in Figure 5a. IGNNet predicted the negative class ($\leq$ 50K) with a narrow margin (0.495). The explanation shows that a single feature (capital-gain=2885) has the highest contribution compared to any other feature value. In the training data, the capital-gain has a maximum value of 99999.0, a minimum value of 0, a mean value of 1068.36, and a 7423.08 standard deviation. To test if the explanation reflects the actual reasoning of IGNNet, we raise the capital-gain value by a smaller value than the standard deviation to be 6000 while leaving the remaining feature values constant, and it turns out to be enough to alter the prediction to a positive ($>$ 50K) with 0.944 as the predicted value. We can also see that the negative score of the capital-gain feature went from -3.67 in the original instance (5a) to -0.82 in the modified instance, as shown in Figure 5b. So the user can adjust the value of an important feature as much as needed to alter the prediction.

The data point, from the Churn dataset, has a positive prediction with a narrow margin (0.565). Consequently, the reduction of the top positively important feature (total day charge), as illustrated in Fig. 6a, may be enough to obtain a negative prediction. The total day charge has a maximum value of 59.76, a minimum of 0.44, a mean of 30.64, and a 9.17 standard deviation in the training set. However, the total day charge is highly correlated with the total day minutes by more than 0.99, as shown in Fig. 7. The high correlation
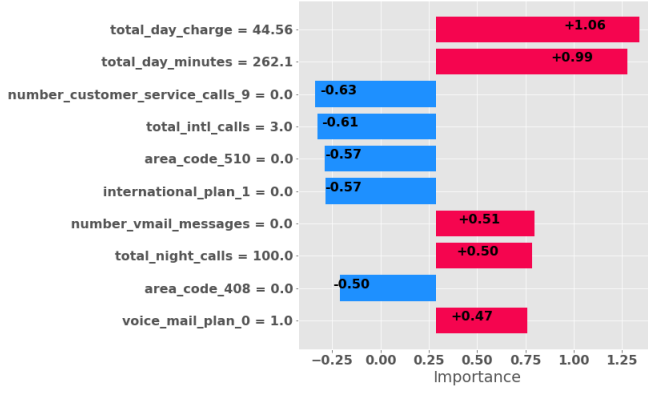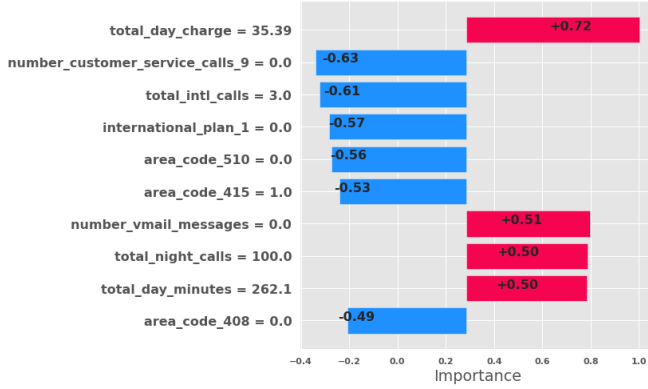


(a) The original data point.



(b) The data point with a modified capital gain value.

Figure 5: Explanation to a single prediction on Adult dataset.

(a) The original data point.



(b) The data point with a modified total day charge value.

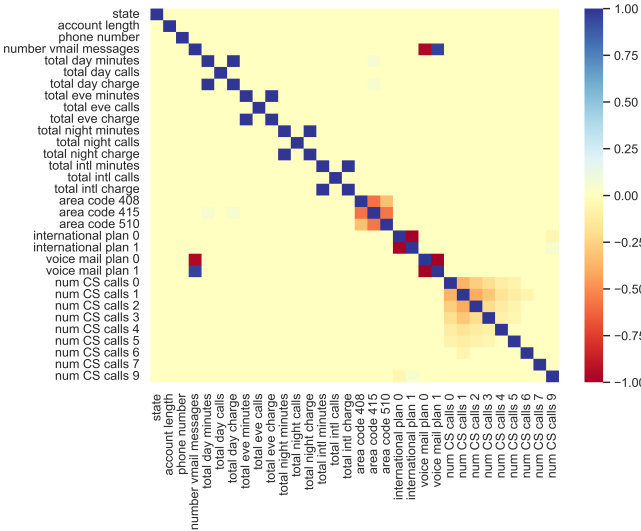Figure 6: The explanation of a single prediction on Churn dataset.



Figure 7: The correlation matrix of the features of the Churn dataset.

effect is obvious in the outcome when the total day charge value is reduced by one standard deviation, from 44.56 to 35.39, as the scores of both total day charge and total day minutes drop from 1.06 and 0.99 to 0.72 and 0.5, respectively, as shown in Fig. 6b. Moreover, the sum of feature scores and bias falls below 0.5 after the sigmoid function, resulting in a negative prediction with a predicted value of 0.297.

## Ablation Study

In this section, we conduct three sets of experiments to determine the effects of different hyperparameter choices on predictive performance. In the first set, we study the effects of changing the weight on the edge from a feature to itself in the adjacency matrix (self-loop). In the second set, we inspect the effect of enforcing a threshold on the correlation value between two features to be included as an edge weight. Finally, the third set of experiments determines the effect of using a different number of message-passing layers.

### Weight on Self-Loop

In the experimental setup, we evaluated IGNNet using a value that forms 70% up to 90% of the weighted summation on the self-loop. In this set of experiments, we evaluate IGNNet using self-loops that form more than 90% of the weighted summation, then we test with a fixed value of 1 for the self-loop, and finally, we evaluate with 0 self-loops. The detailed results are available in Table 2. The null hypothesis that there is no difference in the predictive performance between IGNNet in the default settings, and IGNNet with the three new variations of the self-loops, has been tested using the Friedman test followed by the pairwise Nemenyi tests. The null hypothesis can be rejected at 0.05 level for the difference between the default settings and the tested variations of self-loops, as the default settings showed significantly better performance. The result of the post hoc tests are summarized in Fig. 8.

### Threshold on Correlation Values

In the experimental setup, we set a threshold for a correlation value to be added as an edge between two features, which is 0.2 in general and 0.05 in the case of small correlation values. In this set of experiments, we compare the performance of IGNNet with the abovementioned thresholds and without any thresholds. To test the null hypothesis that there is no difference in the predictive performance, as measured by AUC, between the two variations, and since we compare two models, the Wilcoxon signed-rank test (Wilcoxon 1945) is used here. The null hypothesis may be rejected at the 0.05 level, showing that setting a threshold for the edge weight helps with improving the predictive performance of IGNNet. The results are available in Table 2.

### The Number of Message-Passing Layers

In this final set of experiments, we investigate the effect of reducing the number of the message-passing layers and the total number of parameters in the IGNNet on the predictive
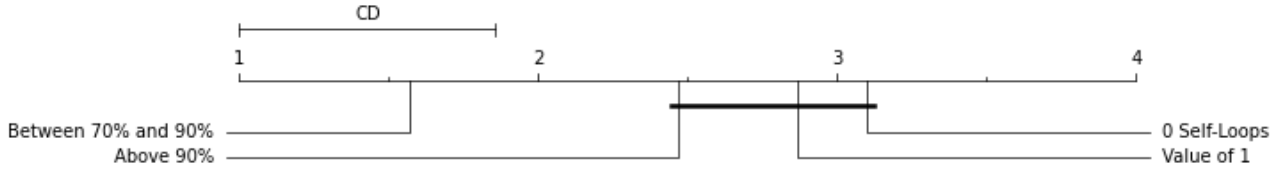
Figure 8: The average rank of IGNNet in the default settings and other IGNNet models with different values for the self-loops on the 30 datasets. The ranking has been done for the predictive performance measured in the AUC.

performance. We compare IGNNet (with the default architecture) to a reduced version of three message-passing layers and another version with just one message-passing layer. For the three layers version, we use the architecture of IGN-Net as mentioned in the experimental setup, but we keep the first three message-passing layers up to the first batch normalization, then we remove the following three message-passing layers and their linear transformations before the final fully-connected feedforward network. The model with one message-passing layer has only the first layer with the subsequent transformations, then all the subsequent layers are removed up until the fully-connected feedforward network is.

Similar to what has been done in the previous subsections, we test the null hypothesis of no difference in the predictive performance between the three architectures using the Friedman test, followed by pairwise post hoc Nemenyi tests. The first test rejects the null hypothesis, and the result of the Nemenyi tests are outlined in Fig. 9. The results in Table 2 demonstrate that the predictive performance generally improves with more layers of message-passing and linear transformations.

## Transparency Evaluation

In this section, we demonstrate the detailed results of the explanations evaluation experiment using 35 datasets. The results show a general trend across the 35 datasets where the explanations obtained using KernelSHAP converge to more similar values to the feature scores produced by IGNNet, given more sampled data. The results are displayed in Figure 10, Figure 11, and Figure 12.

## Information about the Used Datasets and Specifications of the Hardware

This subsection provides a summary of the datasets utilized in the experiments. In Table 3, we provide information about the used datasets, including the number of classes (Num of Classes), the number of features (Num of Features), the size of the dataset, the size of the training, validation, and test splits, as well as, the used correlation threshold of each dataset, the weight on the self-loop (Self-Loop Wt.) and finally the ID of each dataset on OpenML.

The experiments have been performed in a Python environment on an Intel(R) Core(TM) i9-10885H CPU @ 2.40GHz system with 64.0 GB of RAM, and the GPUs are NVIDIA® GeForce® GTX 1650 Ti with 4 GB GDDR6, and

NVIDIA® GeForce® GTX 1080 Ti with 8 GB. All the software and package dependencies are documented with the source code.

Table 2: The results of the ablation study that compares the predictive performance of IGNNet in the default settings to IGNNet with different weights on the self-loop, IGNNet without a threshold on the correlation value to be included as an edge in the adjacency matrix, and IGNNet with different numbers of message-passing layers. The best-performing model is colored in blue, and the second best-performing is colored in light blue.

| Dataset | Default Settings | Self-Loop Experiments | | | Without Correlation Threshold | Number of Layers Experiments | |
|---|---|---|---|---|---|---|---|
| | | Above 90% | Self-Loop Value = 1 | Self-Loop Value = 0 | | 3 Layers | 1 Layer |
| Abalone | 0.88 | 0.88 | 0.844 | 0.842 | 0.88 | 0.874 | 0.848 |
| Ada Prior | 0.905 | 0.905 | 0.891 | 0.88 | 0.896 | 0.904 | 0.903 |
| Adult | 0.917 | 0.914 | 0.916 | 0.88 | 0.916 | 0.917 | 0.911 |
| Bank 32 nh | 0.887 | 0.889 | 0.854 | 0.815 | 0.886 | 0.877 | 0.881 |
| Breast Cancer | 0.999 | 0.999 | 0.998 | 0.998 | 0.999 | 0.998 | 0.998 |
| Churn | 0.88 | 0.865 | 0.871 | 0.9 | 0.905 | 0.872 | 0.868 |
| Credit Card Fraud | 0.987 | 0.974 | 0.942 | 0.913 | 0.952 | 0.938 | 0.937 |
| Delta Ailerons | 0.977 | 0.977 | 0.977 | 0.976 | 0.977 | 0.976 | 0.977 |
| Delta Elevators | 0.951 | 0.952 | 0.946 | 0.946 | 0.951 | 0.95 | 0.949 |
| Electricity | 0.898 | 0.871 | 0.881 | 0.844 | 0.866 | 0.845 | 0.838 |
| Elevators | 0.951 | 0.946 | 0.923 | 0.928 | 0.931 | 0.896 | 0.847 |
| Higgs | 0.762 | 0.791 | 0.696 | 0.707 | 0.758 | 0.727 | 0.717 |
| JM1 | 0.739 | 0.729 | 0.702 | 0.71 | 0.734 | 0.721 | 0.723 |
| Madelon | 0.906 | 0.717 | 0.91 | 0.895 | 0.855 | 0.807 | 0.686 |
| Magic Telescope | 0.907 | 0.929 | 0.89 | 0.91 | 0.919 | 0.901 | 0.873 |
| Mozilla4 | 0.954 | 0.953 | 0.954 | 0.965 | 0.963 | 0.942 | 0.862 |
| MC1 | 0.957 | 0.909 | 0.953 | 0.935 | 0.933 | 0.904 | 0.961 |
| Numerai28.6 | 0.526 | 0.52 | 0.52 | 0.521 | 0.52 | 0.525 | 0.525 |
| PC2 | 0.881 | 0.883 | 0.865 | 0.859 | 0.827 | 0.849 | 0.875 |
| Phishing | 0.99 | 0.987 | 0.989 | 0.993 | 0.992 | 0.988 | 0.987 |
| Phonemes | 0.922 | 0.919 | 0.924 | 0.89 | 0.941 | 0.915 | 0.895 |
| Pollen | 0.525 | 0.491 | 0.516 | 0.524 | 0.483 | 0.489 | 0.498 |
| Satellite | 0.998 | 0.996 | 0.928 | 0.931 | 0.984 | 0.996 | 0.995 |
| Scene | 0.994 | 0.959 | 0.969 | 0.932 | 0.992 | 0.957 | 0.941 |
| Spambase | 0.984 | 0.98 | 0.974 | 0.963 | 0.983 | 0.985 | 0.982 |
| Speed Dating | 0.853 | 0.844 | 0.825 | 0.812 | 0.824 | 0.82 | 0.814 |
| Telco Customer Churn | 0.858 | 0.86 | 0.841 | 0.844 | 0.858 | 0.852 | 0.853 |
| Tic Tac Toe | 0.846 | 0.828 | 0.841 | 0.864 | 0.876 | 0.831 | 0.803 |
| Vehicle sensIT | 0.918 | 0.916 | 0.915 | 0.91 | 0.915 | 0.915 | 0.914 |
| Waveform-5000 | 0.965 | 0.962 | 0.966 | 0.964 | 0.961 | 0.96 | 0.96 |

Table 3: The dataset information.

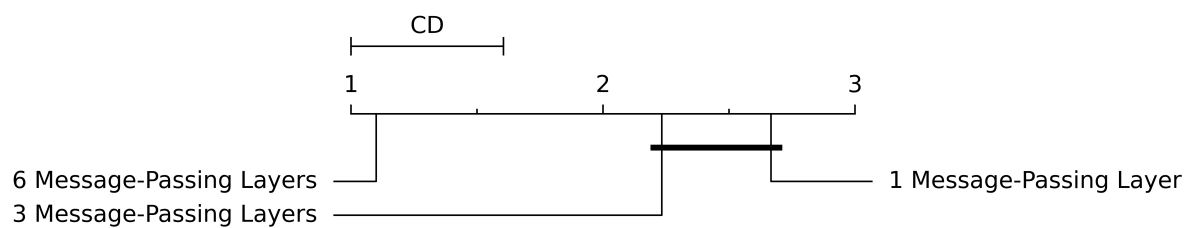| Dataset | Classes | Features | Dataset Size | Training Set | Dev. Set | Test Set | Corr. Threshold | Self-Loop Wt. | Training Epochs | OpenML ID |
|---|---|---|---|---|---|---|---|---|---|---|
| Abalone | 2 | 8 | 4177 | 2506 | 836 | 835 | 0.2 | 20 | 220 | 720 |
| Ada Prior | 2 | 14 | 4562 | 2737 | 913 | 912 | 0.2 | 4 | 216 | 1037 |
| Adult | 2 | 14 | 48842 | 43957 | 2443 | 2442 | 0.2 | 4 | 297 | 1590 |
| Bank 32 nh | 2 | 32 | 8192 | 5734 | 1229 | 1229 | 0.2 | 2 | 540 | 833 |
| Covertype | 7 | 54 | 581012 | 524362 | 27599 | 29051 | 0.2 | 10 | 300 | 1596 |
| Credit Card Fraud | 2 | 30 | 284807 | 270566 | 7121 | 7120 | 0.2 | 3 | 29 | 42175 |
| Delta Ailerons | 2 | 5 | 7129 | 3564 | 1783 | 1782 | 0.2 | 2 | 380 | 803 |
| Electricity | 2 | 8 | 45312 | 36249 | 4532 | 4531 | 0.2 | 3 | 396 | 151 |
| Elevators | 2 | 18 | 16599 | 11619 | 2490 | 2490 | 0.2 | 4 | 353 | 846 |
| HPC Job Scheduling | 4 | 7 | 4331 | 2598 | 867 | 866 | 0.05 | 10 | 435 | 43925 |
| Fars | 8 | 29 | 100968 | 80774 | 10097 | 10097 | 0.05 | 10 | 301 | 40672 |
| 1st Order Theorem Prov. | 6 | 51 | 6118 | 3915 | 979 | 1224 | 0.2 | 30 | 848 | 1475 |
| Helena | 100 | 27 | 65196 | 41724 | 10432 | 13040 | 0.2 | 10 | 550 | 41169 |
| Heloc | 2 | 22 | 10000 | 7500 | 1250 | 1250 | 0.2 | 20 | 234 | 45023 |
| Higgs | 2 | 28 | 98050 | 88245 | 4903 | 4902 | 0.05 | 4 | 394 | 23512 |
| Indian Pines | 8 | 220 | 9144 | 5852 | 1463 | 1829 | 0.2 | 400 | 394 | 41972 |
| Jannis | 4 | 54 | 83733 | 53588 | 13398 | 16747 | 0.05 | 20 | 300 | 41168 |
| JM1 | 2 | 21 | 10885 | 8708 | 1089 | 1088 | 0.2 | 50 | 187 | 1053 |
| LHC Identify Jets | 5 | 16 | 830000 | 749075 | 39425 | 41500 | 0.2 | 10 | 300 | 42468 |
| Madelon | 2 | 500 | 2600 | 1560 | 520 | 520 | 0.05 | 4 | 199 | 1485 |
| Magic Telescope | 2 | 10 | 19020 | 15216 | 1902 | 1902 | 0.2 | 4 | 397 | 1120 |
| MC1 | 2 | 38 | 9466 | 7478 | 994 | 994 | 0.2 | 80 | 198 | 1056 |
| Mozilla4 | 2 | 5 | 15545 | 12436 | 1555 | 1554 | 0.2 | 2 | 280 | 1046 |
| Microaggregation2 | 5 | 20 | 20000 | 12800 | 3200 | 4000 | 0.2 | 15 | 599 | 41671 |
| Numerai28.6 | 2 | 21 | 96320 | 86688 | 4816 | 4816 | 0.2 | 20 | 36 | 23517 |
| Otto Group Product | 9 | 93 | 61878 | 39601 | 9901 | 12376 | 0.05 | 30 | 270 | 45548 |
| PC2 | 2 | 36 | 5589 | 3353 | 1118 | 1118 | 0.2 | 60 | 37 | 1069 |
| Phonemes | 2 | 5 | 5404 | 3782 | 811 | 811 | 0.2 | 1 | 800 | 1489 |
| Pollen | 2 | 5 | 3848 | 2308 | 770 | 770 | 0.2 | 4 | 351 | 871 |
| Satellite | 2 | 36 | 5100 | 2805 | 1148 | 1147 | 0.2 | 60 | 287 | 40900 |
| Scene | 2 | 299 | 2407 | 1203 | 602 | 602 | 0.2 | 40 | 86 | 312 |
| Speed Dating | 2 | 120 | 8378 | 5864 | 1257 | 1257 | 0.2 | 10 | 58 | 40536 |
| Telco Customer Churn | 2 | 19 | 7043 | 4930 | 1057 | 1056 | 0.2 | 3 | 116 | 42178 |
| Vehicle sensIT | 2 | 100 | 98528 | 88675 | 4927 | 4926 | 0.2 | 10 | 172 | 357 |
| waveform-5000 | 2 | 40 | 5000 | 3000 | 1000 | 1000 | 0.2 | 4 | 97 | 979 |

Figure 9: The average rank of IGNNet with 6, 3, and 1 message-passing layers on the 30 datasets. The ranking has been done for the predictive performance measured in the AUC.
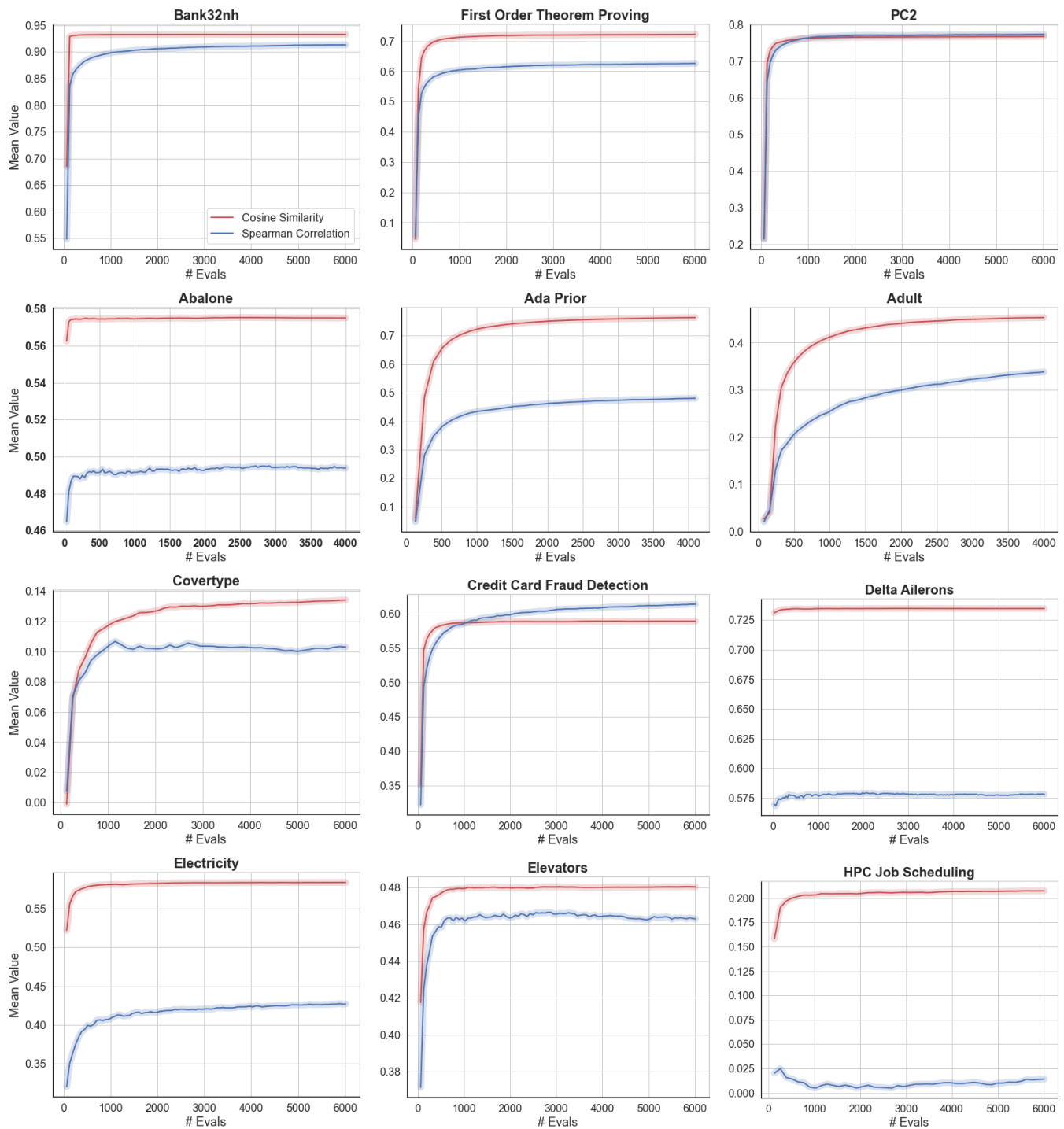
Figure 10: **Comparison of KernelSHAP's approximations and the importance scores obtained from IGNNet.** We measure the similarity of KernelSHAP's approximations to the scores of IGNNet at each iteration of data sampling and evaluation of KernelSHAP. KernelSHAP exhibits improvement in approximating the scores derived from IGNNet with more data sampling.
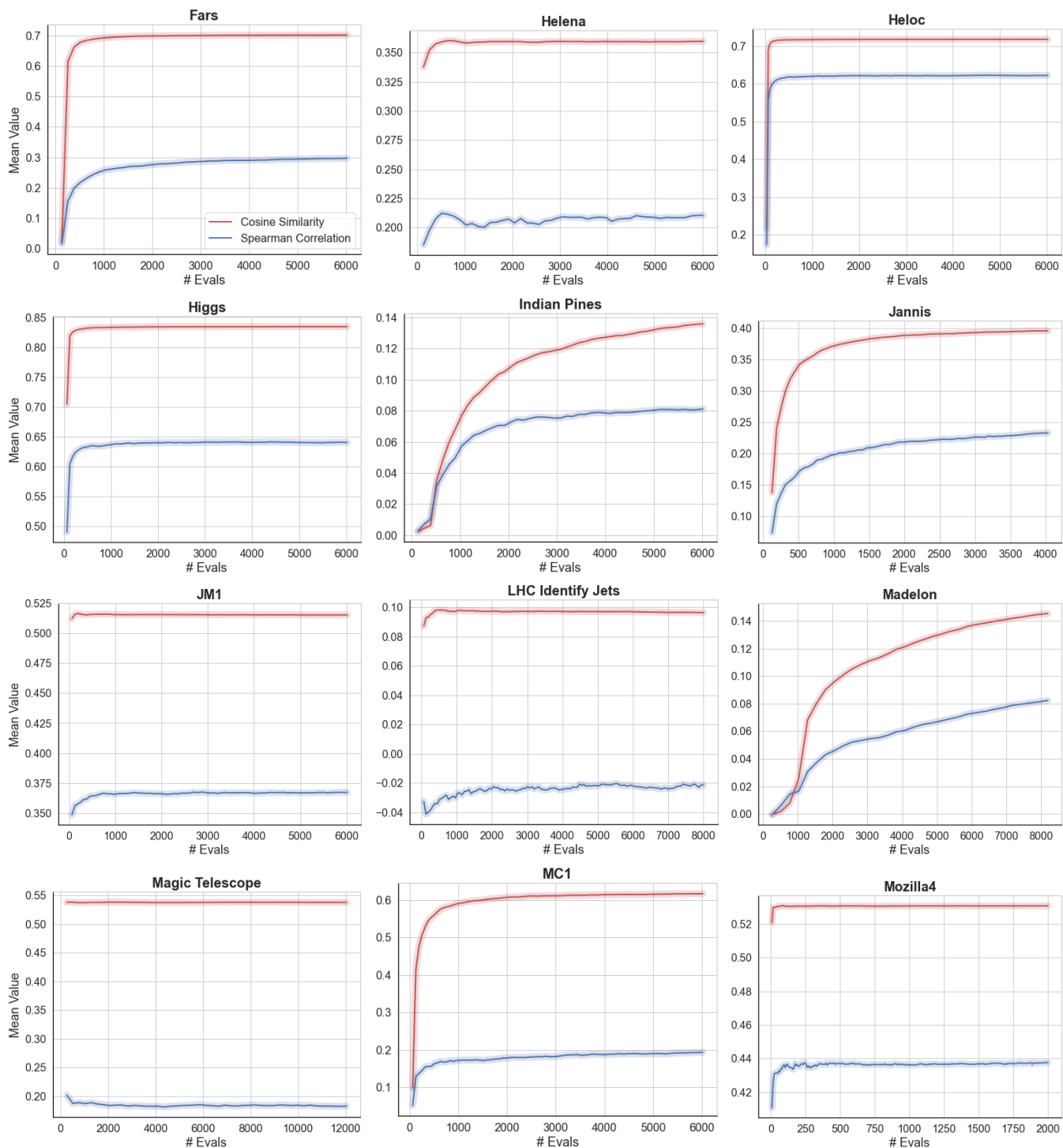
Figure 11: **Comparison of KernelSHAP's approximations and the importance scores obtained from IGNNet.** We measure the similarity of KernelSHAP's approximations to the scores of IGNNet at each iteration of data sampling and evaluation of KernelSHAP. KernelSHAP exhibits improvement in approximating the scores derived from IGNNet with more data sampling.
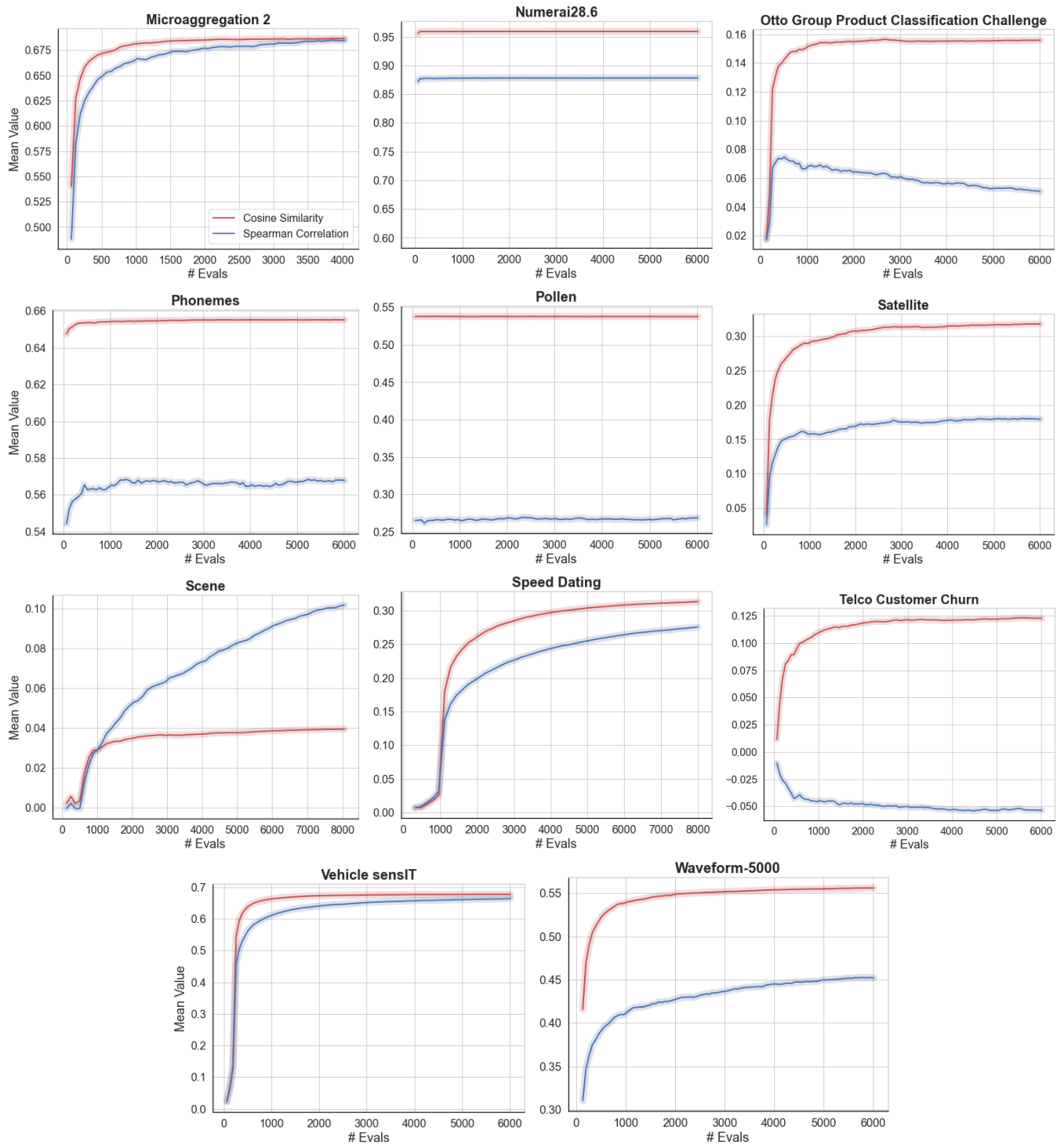
Figure 12: **Comparison of KernelSHAP's approximations and the importance scores obtained from IGNNet.** We measure the similarity of KernelSHAP's approximations to the scores of IGNNet at each iteration of data sampling and evaluation of KernelSHAP. KernelSHAP exhibits improvement in approximating the scores derived from IGNNet with more data sampling.