An Explainer for Temporal Graph Neural Networks

1st Wenchong He* *University of Florida* Gainesville, USA. whe2@ufl.edu 1st Minh N. Vu* University of Florida Gainesville, USA. minhyu@ufl.edu 3rd Zhe Jiang University of Florida Gainesville, USA. zhe.jiang@ufl.edu 4th My T. Thai University of Florida Gainesville, USA. mythai@cise.ufl.edu

Abstract—Temporal graph neural networks (TGNNs) have been widely used for modeling time-evolving graph-related tasks due to their ability to capture both graph topology dependency and non-linear temporal dynamic. The explanation of TGNNs is of vital importance for a transparent and trustworthy model. However, the complex topology structure and temporal dependency make explaining TGNN models very challenging. In this paper, we propose a novel explainer framework for TGNN models. Given a time series on a graph to be explained, the framework can identify dominant explanations in the form of a probabilistic graphical model in a time period. Case studies on the transportation domain demonstrate that the proposed approach can discover dynamic dependency structures in a road network for a time period.

Index Terms—TGNN, graph explanations, interpretable DL

I. Introduction

In the last several years, Graph Neural Networks (GNNs) have been increasingly popular due to their ability to capture the complex relationship and interactions in a system [1]–[9]. GNNs assume static graph data and operate by aggregating information in the local neighborhood of a node. However, many real-world systems are dynamic and evolving over time. For example, in the transportation domain, accurate traffic forecasting requires both temporal features (periodicity and trend) and spatial dependency modeling (the topology of a road network). A Temporal Graph Neural Network (TGNN) combines both a GNN and a recurrent neural network (RNN) to capture both the spatial dependency and temporal dynamics in a system [10]–[12].

Although TGNNs have achieved wide success in modeling spatial networks with temporal dynamics, it is unclear why a model makes certain predictions. Such explanation is crucial because: (1) It improves the transparency and consequently the trust of a TGNN model in a growing number of safety-related applications when being deployed into the real world; (2) It allows end-users to identify interesting dynamic dependency structures in the system. For example, for traffic forecasting, the spatial dependency structure can change due to the evolving characteristics of the traffic flow. In a non-congestion period, the traffic status at upstream roads impacts the traffic status at downstream roads through the transfer effect, and in a congestion period, the traffic status at downstream roads

*Authors contribute equally

impacts the traffic status at upstream roads through the feedback effect [10], [13]. Such dependency dynamics are crucial for a transportation domain expert to identify the influential neighboring roads at different time periods when predicting the traffic flow or speed for one road segment.

The problem of explaining TGNNs is challenging for several reasons. First, to generate an accurate explanation for the temporal predictions, we need to consider the network structural and temporal dependency simultaneously, which increases the problem complexity compared with the static GNN explanations. Second, the explanations identified in some time steps may be insignificant or redundant, it is non-trivial to discover the dominant interesting dependency structure in a time period and eliminate the redundant pattern. Third, the computation complexity is high considering the high complexity of the interpretable domain and quadratic potential temporal windows for explanations.

We propose a novel explainer framework for the temporal graph neural network interpretations based on the static graph explainer. Specifically, the main contributions are:

- To reduce the interpretation complexity, we reformulate the temporal graph neural network model such that the explanations at different temporal snapshots can be done independently to leverage the static graph explainer.
- We propose a pruning approach to discover the temporal dominant interesting explanations from the independent explanations at each time step and eliminate the insignificant and redundant explanations. The proposed pruning algorithms can reduce the temporal search space and improve efficiency.
- The case studies of the TGNN explanations on the transportation domain show that our proposed explainer framework can provide interpretable explanations and discover the dynamic dependency structure of the spatial network.

II. BACKGROUND ON GRAPH NEURAL NETWORK EXPLANATION

This section provides some related background and preliminaries of our paper. The following provides some highlights on GNN and TGNN, the current landscape of explaining predictions made by GNNs, and the PGM-Explainer, which serves as an important component of our method.

GNNs and TGNNs: With the increasing availability of modern graph data and the popularity of graph-related tasks, many Graph-based Neural Networks have been introduced in recent years, such as ChebNets [1], Graph Convolutional

Networks [2], GraphSage [3], Graph Attention Networks [4], among others [5]–[9], [14]. At a high level, these models exploit the non-linear transformation in conventional neural networks to transform input features and aggregate (or propagate) those features via a graph's connectivity. To consider the time-varying features of graph-related tasks, TGNNs combine GNN and RNN models together, e.g., Temporal Graph Convolutional Neural Network (T-GCN) [10], Spatial-Temporal Graph Social Network (STGSN) [12]. A TGNN integrates the GNN architecture (to model the structural topology dependency) and an RNN model (to capture the nonlinear temporal dependency for time series). Recent studies have put significant attention into the TGNN model, such as applications in traffic forecasting, social network analysis, and human trajectory prediction [10], [12], [15].

Explanation methods for GNNs: Similar to some previously studied deep neural network architectures, the GNNs are known to be black-boxes, i.e. it is unclear how they generate predictions. To address this issue, several explanation methods, called explainers, for GNNs have been introduced recently [16]–[18]. The main objective of the explainers is to identify some network components such as nodes, edges, or sub-graphs, that explain, clarify or contribute the most to the model's predictions. The scope of explanations can also vary: some methods focus on explaining the overall behaviors of the model, while others aim to explain the model's prediction on a specific input instance. In this work, we focus on the latter, which is known as the "local explanation" method. Local explanation methods can also be classified based on their methodology: gradients or features-based, perturbationbased, decomposition, and surrogate methods.

Explaining the GNN using PGM-Explainer: This paper aims to extend PGM-Explainer [17], a local explanation method based on a probabilistic graphical model, from GNN explanation to TGNN explanation. The main reason why we choose PGM-Explainer as a base framework is that it can capture the graph dependencies within TGNN's variables. This makes the explainer a strong candidate to explain models with high temporal-spatial dependencies as TGNNs. In the next paragraphs, we provide a more formal description of how PGM-Explainer explains the static GNNs.

We consider the forwarding of a GNN or TGNN to be explained as a function Φ . In practice, since the output of Φ can contain many predictions (for example as in node classification), we denote o the target prediction to be explained. As such, the prediction to be explained can be written as the function $\Phi_o: \mathcal{G} \to \mathcal{K}$, where \mathcal{G} is the set of input graphs and \mathcal{K} is the set of prediction's outputs on that target.

PGM-Explainer represents the GNN's variables via a directed acyclic probabilistic graph, i.e. a Bayesian network [19]. Given a target prediction to be explained o, PGM-Explainer solves the optimal Bayesian network \mathcal{B}^* :

$$\arg \max_{\mathcal{B} \in \mathcal{E}} R_{\Phi,o}(\mathcal{B}), \quad \text{s.t. } |\mathcal{V}(\mathcal{B})| \le M, o \in \mathcal{V}(\mathcal{B}), \quad (1)$$

where \mathcal{E} is the set of all Bayesian networks and o is the random variable corresponding to the target prediction o. The set of

random variables in Bayesian network $\mathcal{V}(\mathcal{B})$ represents the set of explanatory features in the input graph of the model, which is typically a subset of nodes in the input graph. The objective $R_{\Phi,o}(\mathcal{B})$ measure the fitness of the Bayesian network \mathcal{B} with a set of perturbation data generated from the function Φ_o . The constraints in the optimization is to encourage compact solutions and guarantee the target prediction is included in the explanations. One choice of the objective $R_{\Phi,o}(\mathcal{B})$ is the *BIC score*, which is given as follows:

$$R_{\Phi,o}(\mathcal{B}) = \text{score}_{BIC}(\mathcal{B}: \mathcal{D}_o) = l(\hat{\theta}_{\mathcal{B}}: \mathcal{D}_o) - \frac{\log n}{2} \text{Dim}[\mathcal{B}]$$

where \mathcal{D}_o is the perturbation data generated from Φ_o and $\text{Dim}[\mathcal{B}]$ is the dimension of model \mathcal{B} . $\theta_{\mathcal{B}}$ are the parameters of \mathcal{B} and function $l(\theta_{\mathcal{B}}:\mathcal{D}_o)$ is the log-likelihood between the data \mathcal{D}_o and $\theta_{\mathcal{B}}$. $\hat{\theta}_{\mathcal{B}}$ in the score is parameters' value that maximizes the log-likelihood, which is called the maximum likelihood estimator. This choice of objective has been proven to be *consistent* with the data [20].

III. PROBLEM STATEMENT

This section provides a problem formulation for TGNN explanation. Based on our previous introduction of PGM-Explainer in case of static GNNs, we describe the formulation for TGNN and highlight some key properties and challenges of the problem compared to the static case.

The function learnt by a TGNN can be considered as a mapping from a graph G and a sequence of feature matrices to the predictions:

$$Y = \Phi(G; [X_t, \cdots, X_{t+T-1}]),$$
 (2)

where T is the time window of the model. In practice, the total time interval $[t_1,t_2]$ is normally larger than the window T. In that case, the interval $[t_1,t_2]$ is processed into t_2-t_1-T+1 sequences of length T and feed into the TGNN. We consequentially have t_2-t_1-T+1 predictions.

Therefore, the problem of explaining the TGNN's predictions is not only about explaining the function Φ at some specific input sequence $[X_t,\cdots,X_{t+T-1}]$ but also about explaining all predictions during the time interval $[t_1,t_2]$, which includes t_2-t_1-T+1 predictions.

Formally, given a trained TGNN model Φ and the predictions at the temporal interval to be explained, we define our problems as discovering the crucial dependency structures for the model making predictions on the interval. In the followings, we will first address how prediction of each sequence of length T is explained. Then, we describe how all explanations can be combined to explain all predictions on the interval.

IV. APPROACH

Our framework consists of two modules as shown in Figure 1. The first module leverages PGM-explainer to explain TGNN model independently for each time step. Then the second module aims to discovery the dominant interesting explanations from the explanations identified in the first module. In the following we will discuss the two modules and use T-GCN [10] as an example of TGNN model.

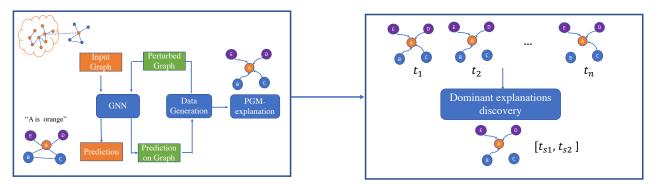


Fig. 1: Our proposed temporal graph neural network explainer framework

A. T-GCN Explanation

A naive way to explain the TGNN Φ on its input $[X_t, \cdots, X_{t+T-1}]$ is to treat the sequence as a single feature matrix and explain the model as in the case of non-temporal models. One drawback of this approach is the temporal information and dependencies among model's variables are not considered and included in the explanations. Furthermore, the number of input's features would be scaled up with a factor of the window T. This introduces significant computation complexity on the explaining process.

To overcome challenges in explaining the TGNN discussed in Sect. III, we reformulate the TGNN equations based on its sequential implementation [10] as follow:

$$Y = \bar{\Phi}(G; X_{t+T-1}; H_{t+T-1}), \tag{3}$$

$$H_{i+1} = \bar{\Phi}(G; X_i; H_i), i = t, \dots, t + T - 2.$$
 (4)

Here, H_i is the hidden features or configurations that the TGNN computes during the sequential forwarding computation. With this formulation, we can consider the problem of explaining the TGNN Φ as the problem of explaining $\bar{\Phi}$ at its last computation with argument H_{t+T-1} . On one hand, the argument H_{t+T-1} is the result of the aggregation of the input features from the past time steps. On the other hand, it can be considered as the configuration of the last computation where the model makes the final prediction. That configuration dictates how the graph's components work together to generate the model's prediction. Thus, we rewrite (Eq. 3) as $Y = \bar{\Phi}_{H_{t+T-1}}(G; X_{t+T-1})$. Then, for each sequence $[X_t, \cdots, X_{t+T-1}]$ the function $\bar{\Phi}_{H_{t+T-1}}$ can be explained by PGM-Explainer as described in Eq. 1.

Specifically, for each time sequence starting with index t, we generate a perturbation data \mathcal{D}_t of $\bar{\Phi}_{H_{t+T-1}}$. Similar to that in PGM-Explainer, \mathcal{D}_t contains two components. First is a set of random seeds determining which nodes of the input graph G is perturbed. Second is an indicator whether the predictions on each nodes is changed by the perturbation on the nodes. The predictions on the perturbations are computed based on the function $\bar{\Phi}_{H_{t+T-1}}$, in which the value of H_{t+T-1} is computed based on the original input data by forwarding the model using Eq. 4. For each \mathcal{D}_t , we use PGM-Explainer

to solve for an optimal Bayesian network \mathcal{B}_t explaining the model at that snapshot.

Note that the explanation Bayesian network \mathcal{B}_t can be different among different snapshots. In analyzing temporal predictions of TGNN, the network structures that can represent many snapshots are much more favorable. Thus, we propose a Temporal Bayesian Information Criterion (TBIC) score measuring the fitness of a given Bayesian network \mathcal{B} with a temporal data $\mathbf{D} = \{\mathcal{D}_t\}_{t=t_s}^{t_e}$:

$$score_{TBIC}(\mathcal{B}: \mathbf{D}) := \frac{1}{t_e - t_s + 1} \sum_{t=t_s}^{t_e} score_{BIC}(\mathcal{B}: \mathcal{D}_t).$$

To mitigate the impact of variations in the data \mathcal{D}_t at each snapshot and grasp a better intuition on how much the edges in the Bayesian network \mathcal{B} help capture the data, we use a normalized version of the TBIC score, called $F_{\mathbf{D}}(\mathcal{B})$. The normalized TBIC is the different between the TBIC score of \mathcal{B} and the Bayesian network with no edge \mathcal{B}_0 :

$$F_{\mathbf{D}}(\mathcal{B}) := \operatorname{score}_{TBIC}(\mathcal{B} : \mathbf{D}) - \operatorname{score}_{TBIC}(\mathcal{B}_0 : \mathbf{D})$$
 (5)

$$= \frac{1}{t_e - t_s + 1} \sum_{t=t_s}^{t_e} \left(\operatorname{score}_{BIC}(\mathcal{B} : \mathcal{D}_t) - \operatorname{score}_{BIC}(\mathcal{B}_0 : \mathcal{D}_t) \right).$$

We can see that each term in the sum of Eq. 5 captures the gain of including the edges in \mathcal{B}_t to fit each data \mathcal{D}_t .

With the normalized TBIC score $F_{\mathbf{D}}$, we now can define the *interesting Bayesian network* and the *temporal dominant interesting Bayesian network*. Specifically, the Bayesian network \mathcal{B} is an *interesting* Bayesian network on a given temporal window $[t_s,t_e]$ if $F_{\mathbf{D}}(\mathcal{B})>B_{threshold}$, where $B_{threshold}$ is a hyper-parameter as a threshold for selecting the interesting Bayesian network. Furthermore, if the temporal window $[t_s,t_e]$ is not a subset of any other interesting Bayesian network temporal window, then \mathcal{B} is a temporal *dominant interesting* Bayesian network on that temporal window. The definition of "dominance" provides a compact explanation.

B. Discover dominant interesting Bayesian networks

In this section, given the independent explanation from PGM-explainer $\{\mathcal{B}_i\}_{i=t_s}^{t_e}$, and the data in each time step $\{\mathcal{D}_i\}_{i=t_s}^{t_e}$, we aim to discover interesting dominant Bayesian

graph in the temporal window. We first described a brute-force approach for temporal pattern discovery. Then we described our proposed pruning approach based on the antimonocity property of dominant interesting Bayesian network.

- 1) Brute-force approach: The brute-force approach has two phases, namely interesting Bayesian network discovery and dominant interesting Bayesian network discovery. In the first phase, for each candidate network, the algorithm scans every possible temporal window and computes the interest score on the window. There are n^2 possible temporal windows and n candidate Bayesian networks. The computation cost is $O(n^3)$ for the first phase. Then in the second phase, the algorithm aims to find dominant interesting Bayesian network whose temporal window is not a subset of that of any other interesting Bayesian network. It scans every pair of interesting Bayesian network's temporal window and eliminates the one that is not dominant. Due to the space limit, we omit the algorithm of the brute-force approach.
- 2) Proposed pruning approach: The brute-force approach does an exhaustive search on all temporal windows for all Bayesian networks. There are lots of redundant computations in this step and requires the second phase to eliminate the redundant Bayesian network explanation. We propose to optimize the brute-force search by a top-down traversal search [21] to leverage the temporal dominant relationship.

Specifically, given a temporal dominant Bayesian network over the temporal window $[t_1,t_2]$, we can conclude that no other temporal dominant Bayesian graph exists in the subtemporal window $[t_1',t_2'] \subset [t_1,t_2]$. Such property inspires a top-down traversal search. We start to calculate the interest measure of each Bayesian graph over the longest temporal window $[t_1,t_2]$. The algorithm recursively reduces the temporal window size if there is no temporal dominant Bayesian network in the temporal window $[t_1,t_2]$. Otherwise, if it finds temporal dominant Bayesian network over $[t_1,t_2]$, the subset of the temporal window of $[t_1,t_2]$ can be pruned out for other Bayesian networks. The computation of interest measures for all Bayesian network over the subset is eliminated.

Then we introduce the detailed implementation of the algorithm. We first construct a directed acyclic graph (DAG) to represent the dominant relationship between temporal windows. Each node represents one temporal window $[t_i,t_j]$ and has two child nodes that are sub-temporal windows and cover one time step less, $[t_i+1,t_j]$ and $[t_i,t_j-1]$. The algorithm uses a queue to do breadth first search on the temporal window. For each node, we evaluate the interest measure for all Bayesian networks over the temporal window. If the node identifies at least one temporal dominant Bayesian network, all the successors of the node will be pruned, otherwise, the algorithm continues evaluating the successors of the node.

Computational complexity analysis: The top-down search approach can potentially eliminate redundant temporal search. In the best case, the algorithm identifies the longest temporal window $[t_s,t_e]$ as the temporal dominant Bayesian network. The computation for one Bayesian network is O(n) for all Bayesian networks. In the worst case, no dominant Bayesian

graph is identified and the algorithm needs to evaluate all n^2 temporal windows. The computational cost for n candidate Bayesian networks is $O(n^3)$.

noend 1 The pruning approach

Input:

- Candidate Bayesian graph $\{\mathcal{B}_i\}_{i=t_s}^{t_e}$
- Perturbed dataset from the model D
- Interest measure function $F_{\mathbf{D}}$

Output:

• All the Interesting dominant Bayesian graph and its dominant time window

```
Candset \leftarrow \emptyset, Create a empty queue Q
Q.enque([t_s,t_e])

while Q not empty do
\mathcal{W} = [t_1,t_2] \leftarrow Q.deque()
for each \mathcal{B}_i do
Scan the temporal window \mathcal{W} and compute the interest measure of Bayesian network \mathcal{B}
if F_{\mathbf{D}}(\mathcal{B}_i) > B_{threshold} then
Add \{\mathcal{B}_i: [t_1,t_2]\} to the CandSet
Prune all subtemporal window of \mathcal{W}
Break
if \mathcal{W} has sub-temporal window then
\mathcal{W}_1 = [t_1+1,t_2], \mathcal{W}_2 = [t_1,t_2-1]
Q.enque(\mathcal{W}_1), Q.enque(\mathcal{W}_2)
return CandSet
```

V. EXPERIMENT RESULTS

In this section, we aim to evaluate our proposed approach on the temporal graph neural network model and present the interpretation results on the transportation domain dataset.

A. Dataset Description

We use the datasets from transportation domain, which contains the traffic speed information on the road network. The dataset SZ-taxi contains the taxi trajectories of Shenzhen from Jan.1, 2015 to Jan.30, 2015. Following the experiments setup, we select 156 major roads of Lanzhou districts as the study area. Each road segment is represented with one node. The adjacency matrix describes the spatial relationship between roads. The feature matrix describes the speed dynamics overtime on each road. The traffic speed was aggregated every 15 minutes. The dataset was split into 60% for training, 20% for validation and 20% for testing. We selected one day (discretize into 96 temporal steps) to do the interpretation of the prediction on the test dataset.

B. Experiment setup

We train T-GCN model [10] using Adam optimizer and set the learning rate to 0.001, batch size to 64 and training epoch to 3000. The accuracy of T-GCN on the test dataset is 0.71. Then we generate the perturbation dataset using 1000 samples. The perturbation threshold is set to 0.01, and perturbation probability is 0.2. For dominant Bayesian network discovery,

the Bayesian score threshold is 1400 based on the validation dataset. All experiments are done on the high performance cluster using 1 NVIDIA A100 GPU with 80 GB GPU memory.

C. Interpretation of results

We present the interpretation results on some representative nodes to visualize the results. The chosen target nodes have around ten two-hop neighbors. The time series of the target node contains both congestion (low speed) and non-congestion (high speed) time slots in one day to be representative for the traffic condition.

Figure 2(a) shows one target node speed ground-truth and the prediction of T-GCN, as well as that of one neighbor node. From the prediction trend, we can observe, the T-GCN model predicts worse at the peak, because the GCN model uses a smooth filter in the spatial network and capture the spatial relationship by constantly moving the filter. Before applying our method to explain the prediction of T-GCN model, we inspect the dependency dynamic of the model in each snapshot. We first construct a standard Bayesian network, which contains all neighborhood nodes within 2-hop and no edges. The Bayesian score of standard network on each temporal snapshot data is shown in Figure 2(b). We can observe in the time 3 to 4 and 11 to 12, the standard Bayesian score is higher, it implies that the target node and neighbors have a higher probability to be independent of each other. On the contrary, around 20 to 22 the standard Bayesian score is lower, meaning there exists an interesting dependency between the target nodes and neighbors.

Secondly, we apply our approach to explain the predictions of T-GCN on each temporal snapshot during a day and obtain the Bayesian network $\{\mathcal{B}_i\}_{i=t_s}^{t_n}$. Then, we measure the fitness of the identified Bayesian network to the data in other time slots using $score(\mathcal{B}:D_i)$. To evaluate the relative improvement of adding the dependency edges, we use the difference of the Bayesian score between each network \mathcal{B}_i and the standard Bayesian network \mathcal{B}_0 . Figure 3 shows the score results of three dominant interesting Bayesian network. We can observe around the time window 3 to 4, 6 to 8 and 19 to 23, the average score is higher. We show the explanation results of the target node in Figure 4. Three dominant Bayesian network were discovered in the time window [3, 4], [6, 8] and [19, 23]. We can observe that in the early morning period, the influential nodes to the target node is nodes 12, 126, 145, while in the evening period, more dependency structure is observed. Almost all neighbor nodes within 2-hop have high influence on the target node in the evening. Based on the discovery, it is reasonable to hypothesize that the early morning period traffic in this target node is less rush than in the evening.

As a comparison we provide the explanation results for another target node in Figure 5. We can observe for this target node, the dominant Bayesian network discovered around [8,10] (Figure 5 (b)) has much more dependency structure than the dominant Bayesian network in [21,23] (Figure 5 (c)), so we hypothesize that for this target node, the time window [8,10] maybe rush-hour period.

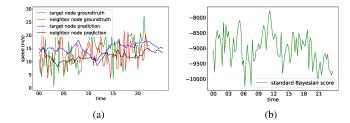


Fig. 2: Ground-truth and prediction of one target node and neighbor (a), standard Bayesian score of the target node (b) (The x axis represents the hour in one day).

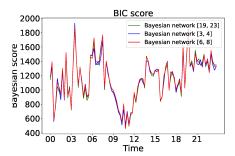


Fig. 3: The Bayesian score evaluated at the day for three candidate Bayesian graph

Analysis of Bayesian score threshold: The selection of interesting Bayesian network depends on the Bayesian score threshold. A higher threshold can find more critical patterns of the network dependency structure, but the related dominant windows are shorter. On the other hand, a lower threshold can find dominant Bayesian network with a longer time window but there may be redundancy. The selection of the threshold can be based on the above trade-off and domain knowledge.

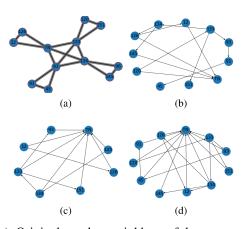


Fig. 4: (a). Original two-hop neighbors of the target node (TN), (b-d). Interesting dominant Bayesian discovered in temporal window: (b). [3, 4], (c). [6, 8], (d). [19, 23]

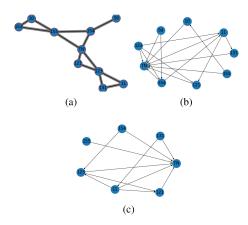


Fig. 5: (a). Original two-hop neighbors of the target node (TN). Interesting dominant Bayesian network discovered in temporal window: (b). [8, 10], (c). [21, 23].

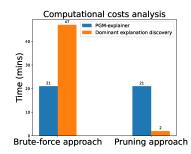


Fig. 6: Computational cost analysis

D. Analysis of Computational costs

Then we analyze the computational costs for brute-force and pruning approach. As Figure 6 shows, the PGM-explainer module requires the same time for two approaches (22 mins). For dominant explanation discovery module, our pruning approach improve significantly than the brute-force approach.

VI. CONCLUSION

In this paper, we propose a novel explainer framework for TGNN model. The framework consists of two modules that first explain each time slot prediction independently and then discovery the dominant interesting explanations in a time period. Case studies show that the proposed approach can discover the dynamic structure dependency in TGNN model.

ACKNOWLEDGEMENT

This material is based upon work supported by the National Science Foundation (NSF) under Grant No. FAI-1939725, SCH-2123809, IIS-2147908, IIS-2207072, and OAC-2152085.

REFERENCES

[1] M. Defferrard, X. Bresson, and P. Vandergheynst, "Convolutional neural networks on graphs with fast localized spectral filtering," in *Proceedings of the 30th International Conference on Neural Information Processing Systems*, 2016, p. 3844–3852.

- [2] T. N. Kipf and M. Welling, "Semi-Supervised Classification with Graph Convolutional Networks," in *Proceedings of the 5th International Con*ference on Learning Representations, 2017.
 [3] W. Hamilton, Z. Ying, and J. Leskovec, "Inductive representation
- [3] W. Hamilton, Z. Ying, and J. Leskovec, "Inductive representation learning on large graphs," in *Advances in Neural Information Processing Systems* 30, 2017, pp. 1024–1034.
- [4] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio, "Graph Attention Networks," *International Conference on Learning Representations*, 2018.
- [5] R. Levie, F. Monti, X. Bresson, and M. M. Bronstein, "Cayleynets: Graph convolutional neural networks with complex rational spectral filters," *IEEE Transactions on Signal Processing*, vol. 67, no. 1, pp. 97–109, Jan 2019.
- [6] F. Monti, K. Otness, and M. M. Bronstein, "Motifnet: A motif-based graph convolutional network for directed graphs," 2018 IEEE Data Science Workshop (DSW), pp. 225–228, 2018.
- [7] K. Xu, W. Hu, J. Leskovec, and S. Jegelka, "How powerful are graph neural networks?" in *International Conference on Learning Represen*tations, 2019.
- [8] Z. Ying, J. You, C. Morris, X. Ren, W. Hamilton, and J. Leskovec, "Hierarchical graph representation learning with differentiable pooling," in *Advances in Neural Information Processing Systems* 31, 2018, pp. 4800–4810.
- [9] Z. Xinyi and L. Chen, "Capsule graph neural network," in *International Conference on Learning Representations*, 2019.
- [10] L. Zhao, Y. Song, C. Zhang, Y. Liu, P. Wang, T. Lin, M. Deng, and H. Li, "T-gcn: A temporal graph convolutional network for traffic prediction," *IEEE Transactions on Intelligent Transportation Systems*, vol. 21, no. 9, pp. 3848–3858, 2019.
- [11] X. Wang, Y. Ma, Y. Wang, W. Jin, X. Wang, J. Tang, C. Jia, and J. Yu, "Traffic flow prediction via spatial temporal graph neural network," in Proceedings of The Web Conference 2020, 2020, pp. 1082–1092.
- [12] S. Min, Z. Gao, J. Peng, L. Wang, K. Qin, and B. Fang, "Stgsn—a spatial-temporal graph neural network framework for time-evolving social networks," *Knowledge-Based Systems*, vol. 214, p. 106746, 2021.
- [13] C. J. Dong, C. F. Shao, C.-X. Zhuge, and M. Meng, "Spatial and temporal characteristics for congested traffic on urban expressway," *Journal of Beijing University of Technology*, vol. 38, no. 8, pp. 1242– 1246, 2012.
- [14] J. Lee, I. Lee, and J. Kang, "Self-attention graph pooling," in 36th International Conference on Machine Learning, ICML 2019, 2019, pp. 6661–6670.
- [15] A. Mohamed, K. Qian, M. Elhoseiny, and C. Claudel, "Social-stgcnn: A social spatio-temporal graph convolutional neural network for human trajectory prediction," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 14424–14432.
- [16] Z. Ying, D. Bourgeois, J. You, M. Zitnik, and J. Leskovec, "GN-Nexplainer: Generating explanations for graph neural networks," in *Advances in Neural Information Processing Systems* 32, 2019, pp. 9244–9255.
- [17] M. Vu and M. T. Thai, "Pgm-explainer: Probabilistic graphical model explanations for graph neural networks," in *Advances in Neural Informa*tion Processing Systems, H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, Eds., vol. 33. Curran Associates, Inc., 2020, pp. 12 225–12 235.
- [18] T. Schnake, O. Eberle, J. Lederer, S. Nakajima, K. Schütt, K.-R. Mueller, and G. Montavon, "Higher-order explanations of graph neural networks via relevant walks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. PP, pp. 1–1, 09 2021.
- [19] J. Pearl, "Chapter 3 markov and bayesian networks: Two graphical representations of probabilistic knowledge," in *Probabilistic Reasoning* in *Intelligent Systems*, 1988, pp. 77 – 141.
- [20] D. Koller and N. Friedman, Probabilistic Graphical Models: Principles and Techniques - Adaptive Computation and Machine Learning. The MIT Press, 2009.
- [21] X. Zhou, S. Shekhar, and D. Oliver, "Discovering persistent change windows in spatiotemporal datasets: A summary of results," in *Proceedings* of the 2nd ACM SIGSPATIAL International Workshop on Analytics for Big Geospatial Data, 2013, pp. 37–46.