
TBoost: Gradient Boosting Temporal Graph Neural Networks

Anonymous Author(s)

Affiliation

Address

email

1 Supplementary Material

2 1.1 Algorithm

Algorithm 1: TBoost Training Loop

```
1 Function Training_Loop( $\mathcal{G}, \mathbf{x}_u, \mathbf{x}_v, \mathbf{x}_{u,v}(t), y(t)$ ):
2    $\mathbf{x}_{u,v}^{gbdt}(t) = [\mathbf{x}_u, \mathbf{x}_v, \mathbf{x}_{u,v}(t), TE(t)]$ 
3    $y^{gbdt}(t) = y(t)$ 
4   for iter  $i=1$  to  $N$  do
5     // Train  $k$  decision trees
6      $f_i \leftarrow \arg \min_{f_i} \mathcal{L}^{gbdt}(f_i(\mathbf{x}_{u,v}^{gbdt}(t)), y(t))$ 
7      $f = f + f_i$ ;
8     // Prepare GNN input features
9      $\bar{\mathbf{x}}_{u,v}(t) = [\mathbf{x}_{u,v}^{gbdt}(t), \hat{y}_{u,v}^{gbdt}(t)]$ 
10     $\mathbf{x}_{u,v}^{gnn}(t) = [\mathbf{x}_u, \mathbf{x}_v, \bar{\mathbf{x}}_{u,v}(t)]$ 
11    // Train  $l$  epochs of GNN
12     $\theta, -\eta \frac{\partial \mathcal{L}^{gnn}}{\partial \mathbf{x}_{u,v}^{gnn}(t)} \leftarrow$  Equation 4
13    // update the target for the next iteration of GBDT
14     $y^{gbdt}(t) = -\eta \frac{\partial \mathcal{L}^{gnn}}{\partial \mathbf{x}_{u,v}^{gnn}(t)}$ 
15  end for
16  Output: GBDT  $f$ , GNN  $g_\theta$ 
```

3 1.2 Baselines

4 We evaluate TBoost against the following baselines:

- 5 • **GBDT Baselines:** CatBoost[9] and XGBoost[1] are commonly used gradient boosting
6 algorithms. They differ in handling categorical variables, missing values, and overfitting
7 prevention. The models are trained on the concatenated node and edge features of datasets.
8 We also implement another version of GBDTs where the model is trained on concatenated
9 temporal encoding features as given in Equation 3. These models are called XGBoost+ and
10 CatBoost+ respectively.
- 11 • **Temporal GNN Baselines:** In our study, we utilized several state-of-the-art temporal graph
12 neural network (GNN) models, including JODIE [4], DyRep [11], and TGAT [12], which are
13 specifically designed for learning on dynamic graphs. JODIE [4] is a neural network model
14 that learns bipartite embeddings for both nodes and edges in a dynamic graph by modeling

15 the interactions between nodes as a sequence of edge snapshots over time. DyRep [11],
 16 on the other hand, uses temporal point processes to model the evolution of the graph as a
 17 sequence of events and predict user-item interactions in temporal graphs where the structure
 18 of the graph changes over time. TGAT [12] extends the basic graph attention network (GAT)
 19 to incorporate temporal information, allowing it to capture the dynamics of the graph over
 20 time. TGN [10] combines RNN-based temporal modeling with graph attention convolution.
 21 It captures temporal information using an RNN and incorporates self-attention mechanisms
 22 to model spatial and temporal dynamics jointly. To evaluate the baselines, we used public
 23 PyTorch implementation of the repositories released by the authors.

24 1.3 Datasets

Table 1: Dataset statistics

Dataset	Vertices	Edges	Timestamps	Classes
Wikipedia	9,227	1,57,474	1,52,757	2
Reddit	10,984	6,72,447	6,69,065	2
MOOC	7145	411,749	345,600	2
Fashion	437	3,176	344	5
Luxury	5,400	34,278	2,510	5
Office	129,466	800,357	5,348	5
Payments	367,221	1,063,398	544,724	2
BankSim	4,162	594,643	180	2

25 We consider real-world datasets Wikipedia [4], Reddit [4] and MOOC [4] along with simulated
 26 financial datasets Payments [2], BankSim [5] for the temporal node classification task and Amazon
 27 [7] for the temporal edge classification task. The dataset statistics are shown in Table 1:

- 28 • **Wikipedia** [4] The dataset consists of one month’s worth of Wikipedia page modifications,
 29 focusing on the top 1,000 pages with the most revisions and 8,227 editors with at least five
 30 edits as users. User modifications are transformed into 172-dimensional LIWC [8] feature
 31 vectors.
- 32 • **Reddit** [4] This dataset is a temporal graph of active users and their posts under subreddits
 33 as 11,000 nodes, 700,000 temporal edges, and dynamic labels suggesting whether a user
 34 has been banned from posting. Edge feature vectors are created from the user post in the
 35 temporal setting. Here, we are trying to predict whether a user has been banned to post.
- 36 • **MOOC** [4]: This public dataset consists of actions, e.g., viewing a video, submitting an
 37 answer, etc., done by students on a MOOC online course. This dataset consists of 7,047
 38 users interacting with 98 items (videos, answers, etc.) resulting in over 411,749 interactions.
 39 The task is to predict one of the 4,066 drop-out events.
- 40 • **Payments Data**¹ [2] The data consists of subject-centric transaction records aimed at de-
 41 tecting fraudulent activities, encompassing both normal and abnormal transactions with
 42 predefined probabilities. It was generated using an AI planning execution simulator, translat-
 43 ing output planning traces into a tabular format, with parameters like the number of clients,
 44 time duration, and fraud probabilities configurable.
- 45 • **BankSim** [5] BankSim is an agent-based simulator generating synthetic data for fraud
 46 detection research, based on Spanish bank transaction data. It models merchant-customer
 47 relationships and aims to represent normal payments and known fraud patterns.
- 48 • **Amazon** [Fashion, Luxury and Office] [7] The Amazon review dataset is a vast collection
 49 of product reviews from Amazon.com, comprising reviews with text, product IDs, reviewer
 50 IDs, ratings, and timestamps. Word2Vec [6] embeddings of the reviews are utilized as
 51 features, while user ratings serve as the labels for analysis.

¹This publication includes or references synthetic data provided by J.P. Morgan.

52 1.4 Model Training

53 TBoost is an end-to-end framework that is trained in an alternate fashion using Adam [3] optimizer.
 54 TBoost is capable of accommodating any temporal GNN model along with GBDT. We experimented
 55 with CatBoost as a GBDT module. In temporal GNN, either TGAT or TGN is employed as these
 56 models perform well over other baselines. Based on our experiments with TBoost, we found that
 57 a single layer of GNN with two attention heads is sufficient to model the data for attention-based
 58 models. We used a learning rate of $1e^{-3}$ for Wikipedia and Reddit and a learning rate of $1e^{-5}$
 59 for other datasets. The neighborhood size was set to 10 for all models, with a batch size of 30
 60 for Wikipedia, 200 for Reddit, and 100 for the other datasets. For tree-based models, we used a
 61 learning rate of $1e^{-3}$, a maximum depth of 7, and a maximum of 1000 iterations. We also employed
 62 early stopping with a patience of 5. Trained models are evaluated for node classification and edge
 63 classification tasks. For the node-classification task, due to the label imbalance in the datasets, we
 64 employ the area under the ROC curve (AUC) as done in prior arts [12, 10]. In the multiclass edge
 65 classification task, the number of classes is 5 for all the datasets in our experimentation. We use the
 66 macro-F1 score as the evaluation metric for these tasks. We carry out each experiment 10 times and
 67 report the mean numbers in Table 2 and Table 3 in the main material.

68 1.5 Ablation study

Table 2: Impact of time-encoding on GBDT for the temporal node classification task. In TE₁ temporal feature is not used, in TE₂ absolute value of timestamp is used, in TE₃ Bochner’s kernel is used and in TE₄ inter arrival time is used.

Model	CatBoost				XGBoost			
	TE ₁	TE ₂	TE ₃	TE ₄	TE ₁	TE ₂	TE ₃	TE ₄
Reddit	56.62	56.23	54.30	56.59	53.34	57.30	51.45	54.78
Wiki	59.30	63.84	62.54	65.75	55.50	64.74	59.37	64.91
MOOC	59.73	60.11	60.26	60.00	59.82	60.10	57.75	59.81
Payments	52.70	50.10	52.50	56.10	52.20	53.07	51.11	54.60
BankSim	52.70	51.47	51.30	52.90	52.20	50.30	51.06	52.25
Average	56.21	56.35	56.18	58.27	54.61	56.70	54.10	57.27

69 **Effect of temporal encoding:** GBDT models, while widely used for various supervised learning
 70 tasks, are not inherently designed to handle time-series data. There is no built-in mechanism within
 71 the algorithm to handle temporal aspects. Several approaches have been explored in the literature
 72 to encode time information in GBDTs. These methods include incorporating time as a feature
 73 (TE₂) in Table 2, adopting functional encoding inspired by Bochner’s Kernel (TE₃), and utilizing
 74 inter-event time (TE₄), among others(Refer Table 2). Due to the discrete nature of GBDTs, we
 75 used a non-trainable version of the temporal kernels. We report the ROC-AUC in Table 2 for the
 76 node classification task averaged over 10 runs. We conducted experiments using two extensively
 77 used GBDT algorithms, CatBoost [9] and XGBoost [1]. Our results showed consistent performance
 78 improvements by incorporating time information, but the candidate for best encoding changed across
 79 different datasets. On average, Bochner’s encoding did not yield substantial improvements in GBDTs,
 80 performing similarly to models without time information. On the other hand, adding time as a feature
 81 had better results. Further, the best results were obtained with inter-event time.

82 **Effect of pretraining:** In the temporal GNN literature, one of the common ways to benchmark
 83 node classification is by pre-training the GNN with an unsupervised objective like contrastive link
 84 prediction. For downstream tasks, the GNNs are initialized with the pre-trained weights and are
 85 kept non-trainable. Only the task-specific head is trained in a supervised fashion. In this section, we
 86 tried to understand the impact of pretraining and weight fixing in TBoost for the node classification
 87 task. The results are shown in the Table 3. In the first setting (TBoost₁), we train Tboost without
 88 using pre-trained weights. In the second setting (TBoost₂), we initialize the GNN with pre-trained
 89 weights that are eventually frozen during training. Only the classifier layer is trained through
 90 backpropagation. In the third setting (TBoost₃), we initialize the GNN with pre-trained weights
 91 but don’t freeze them. We keep the entire framework trainable. Table 3 shows that pre-training
 92 helps boost performance in most cases. We have used the same in our implementation of TBoost.

Table 3: Effect of pretraining on link classification. TBoost₁ is without pre-trained GNN model. TBoost₂ is initialized with pre-trained TGAT but the GNN weights are frozen. In TBoost₃, the pre-trained weights are kept trainable.

Model	Reddit	Wikipedia	MOOC	Payments	BankSim
TBoost ₁	65.80	85.28	51.31	56.31	77.65
TBoost ₂	66.03	87.50	59.32	57.37	78.10
TBoost ₃	66.09	88.40	58.40	62.10	78.15

- ⁹³ Keeping the GNN parameters trainable gives a further boost compared to the frozen parameter setting.
⁹⁴ However, for a fair comparison with baselines in Table 3, we report only the fixed weighted results.