

Scalable Graph Neural Networks via Bidirectional Propagation

Ming Chen

Renmin University of China
chenmming@ruc.edu.cn

Zhewei Wei*

Renmin University of China
zhewei@ruc.edu.cn

Bolin Ding

Alibaba Group
bolin.ding@alibaba-inc.com

Yaliang Li

Alibaba Group
yaliang.li@alibaba-inc.com

Ye Yuan

Beijing Institute of Technology
yuan-ye@bit.edu.cn

Xiaoyong Du

Renmin University of China
duyong@ruc.edu.cn

Ji-Rong Wen

Renmin University of China
jrw@ruc.edu.cn

Abstract

Graph Neural Networks (GNN) is an emerging field for learning on non-Euclidean data. Recently, there has been increased interest in designing GNN that scales to large graphs. Most existing methods use "graph sampling" or "layer-wise sampling" techniques to reduce training time. However, these methods still suffer from degrading performance and scalability problems when applying to graphs with billions of edges. This paper presents GBP, a scalable GNN that **utilizes a localized bidirectional propagation process from both the feature vectors and the training/testing nodes**. Theoretical analysis shows that GBP is the first method that achieves sub-linear time complexity for both the precomputation and the training phases. An extensive empirical study demonstrates that GBP achieves state-of-the-art performance with significantly less training/testing time. Most notably, GBP can deliver superior performance on a graph with over 60 million nodes and 1.8 billion edges in less than half an hour on a single machine.

1 Introduction

Recently, the field of Graph Neural Networks (GNNs) has drawn increasing attention due to its wide range of applications such as social analysis [23, 20, 28], biology [10, 26], recommendation systems [36], and computer vision [39, 7, 13]. Graph Convolutional Network (GCN) [15] adopts a message-passing approach and gathers information from the neighbors of each node from the previous layer to form the new representation. The vanilla GCN uses a full-batch training process and stores each node's representation in the GPU memory, which leads to limited scalability. On the other hand, training GCN with mini-batches is difficult, as the neighborhood size could grow exponentially with the number of layers.

These techniques can be broadly divided into three categories: 1) **Layer-wise sampling methods**: GraphSAGE [11] proposes a neighbor-sampling method to sample a fixed number of neighbors for each node. VRGCN [6] leverages historical activations to restrict the number of sampled nodes

*Zhewei Wei is the corresponding author. Work partially done at Gaoling School of Artificial Intelligence, Beijing Key Laboratory of Big Data Management and Analysis Methods, MOE Key Lab DEKE, Renmin University of China, and Pazhou Lab, Guangzhou, 510330, China.

and reduce the variance of sampling. FastGCN [5] samples nodes of each layer independently based on each node’s degree and keeps a constant sample size in all layers to achieve scales linearly. LADIES [40] further proposes a layer-dependent sampler to constrain neighbor dependencies, which guarantees the connectivity of the sampled adjacency matrix. 2) **Graph Sampling methods**: ClusterGCN [8] builds a complete GCN from clusters in each mini-batch. GraphSAINT [37] proposes several light-weight graph samplers and introduces a normalization technique to eliminate biases of mini-batch estimation. 3) **Linear Models**: SGC [30] computes the product of the feature matrix and the k -th power of the normalized adjacency matrix during the preprocessing step and performs standard logistic regression to remove redundant computation. PPRGo [4] uses Personalized PageRank to capture multi-hop neighborhood information and uses a forward push algorithm [2] to accelerate computation.

While the above methods significantly speed up the training time of GNNs, they still suffer from three major **drawbacks**. First of all, **the time complexity is linear to m** , the number of edges in the graph. In theory, this complexity is undesirable for scalable GNNs. Secondly, as we shall see in Section 4, the existing scalable GNNs, such as GraphSAINT, LADIES, and SGC, **fail to achieve satisfying results** in the semi-supervised learning tasks. Finally and most importantly, none of the existing methods can **offer reliable performance on billion-scale graphs**.

Our contributions. In this paper, we first carefully analyze the theoretical complexity of existing scalable GNNs and explain why they cannot scale to graphs with billions of edges. Then, we present GBP (Graph neural network via **Bidirectional Propagation**), a scalable Graph Neural Network with sub-linear time complexity in theory and superior performance in practice. **GBP performs propagation from both the feature vector and the training/testing nodes, yielding an unbiased estimator for each representation.** Each propagation is executed in a localized fashion, leading to sub-linear time complexity. After the bidirectional propagation, each node’s representation is fixed and can be trivially trained with mini-batches. The empirical study demonstrates that GBP consistently improves the performance and scalability across a wide variety of datasets on both semi-supervised and fully-supervised tasks. Finally, we present the **first empirical study on a graph with over 1.8 billion edges**. The result shows that GBP achieves superior results in less than 2,000 seconds on a moderate machine.

2 Theoretical analysis of existing methods

In this section, we review some of the recent scalable GNNs and analyze their theoretical time complexity. We consider an undirected graph $G=(V, E)$, where V and E represent the set of vertices and edges, respectively. For ease of presentation, we assume that **G is a self-looped graph** [15], with a self-loop attached to each node in V . Let $n = |V|$ and $m = |E|$ denote the number of vertices and edges in G , respectively. Each node is associated with an F -dimensional feature vector and we use $\mathbf{X} \in \mathcal{R}^{n \times F}$ to denote the feature matrix. We use \mathbf{A} and \mathbf{D} to represent the adjacency matrix and the diagonal degree matrix of G , respectively. For each node $u \in V$, $N(u)$ is the set of neighbors of u , and $d(u) = |N(u)|$ is the degree of u . We use $d = \frac{m}{n}$ to denote the average degree of G . Following [15], we define the normalized adjacency matrix of G as $\tilde{\mathbf{A}} = \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}$. The $(\ell + 1)$ -th layer $\mathbf{H}^{(\ell+1)}$ of the vanilla GCN is defined as

$$\mathbf{H}^{(\ell+1)} = \sigma \left(\tilde{\mathbf{A}} \mathbf{H}^{(\ell)} \mathbf{W}^{(\ell)} \right), \quad (1)$$

where $\mathbf{W}^{(\ell)}$ is the learnable weight matrix and $\sigma(\cdot)$ is the activation function. The training and inference time complexity of a GCN with L layers can be bounded by $O(LmF + LnF^2)$, where $O(LmF)$ is the total cost of the sparse-dense matrix multiplication $\tilde{\mathbf{A}} \mathbf{H}^{(\ell)}$, and $O(LnF^2)$ is the total cost of the feature transformation by applying $\mathbf{W}^{(\ell)}$. At first glance, $O(LnF^2)$ seems to be the dominating term, as the average degree d on scale-free networks is usually much smaller than the feature dimension F and hence $LnF^2 > LndF = LmF$. However, in reality, **feature transformation can be performed with significantly less cost due to better parallelism of dense-dense matrix multiplications**. Consequently, **$O(LmF)$ is the dominating complexity** term of GCN and performing full neighbor propagation $\tilde{\mathbf{A}} \mathbf{H}^{(\ell)}$ is the main bottleneck for achieving scalability.

In order to speed up GCN training, a few recent methods use various techniques to approximate the full neighbor propagation $\tilde{\mathbf{A}}\mathbf{H}^{(\ell)}$ and enable mini-batch training. We divide these methods into three categories and summarize their time complexity in Table 1.

Layer-wise sampling methods sample a subset of the neighbors at each layer to reduce the neighborhood size. GraphSAGE [11] samples s_n neighbors for each node and only aggregates the embeddings from the sampled nodes. With a batch-size b , the cost of feature propagation is bounded by $O(bs_n^L F)$, and thus the total per-epoch cost of GraphSAGE is $O(ns_n^L F + ns_n^{L-1} F^2)$. This complexity grows exponentially with the number of layers L and is not scalable on large graphs. Another work [6] based on node-wise sampling further reduces sampling variance to achieve a better convergence rate. However, it suffers from worse time and space complexity. FastGCN [5] and LADIES [40] restrict the same sample size across all layers to limit the exponential expansion. If we use s_l to denote the number of nodes sampled per layer, the per-batch feature propagation time is bounded by $O(Ls_l^2 F)$. Since $\frac{n}{s_l}$ batches are needed in an epoch, it follows that the per-epoch forward propagation time is bounded by $O(Lns_l F + LnF^2)$. Mini-batch training significantly accelerates the training process of the layer-wise sampling method. However, the training time complexity is still linear to m as the number of samples s_l is usually much larger than the average degree d . Furthermore, it has been observed in [8] that the overlapping nodes in different batches will lead to high computational redundancy, especially in fully-supervised learning.

Graph sampling methods sample a sub-graph at the beginning of each batch and perform forward propagation on the same subgraph across all layers. Cluster-GCN [8] uses graph clustering techniques [14] to partition the original graph into several sub-graphs, and samples one sub-graph to perform feature propagation in each mini-batch. In the worst case, the number of clusters in the graph is 1, and Cluster-GCN essentially becomes vanilla GCN in terms of time complexity. GraphSAINT [37] samples a certain amount of nodes and uses the induced sub-graph to perform feature propagation in each mini-batch. Let b denote the number of sampled node per-batch, and $\frac{n}{b}$ denote the number of batches. Given a sampled node u , the probability that a neighbor of u is also sampled is b/n . Therefore, the expected number of edges in the sub-graph is bounded by $O(b^2 d/n)$. Summing over $\frac{n}{b}$ batches follows that the per-epoch feature propagation time of GraphSAINT is bounded by $O(LbdF)$, which is sub-linear to the number of edges in the graph. However, GraphSAINT requires a full forward propagation in the inference phase, leading to the $O(LmF + LnF^2)$ time complexity.

Linear model removes the non-linearity between each layer in the forward propagation, which allows precomputation of the final feature propagation matrix and result in an optimal training time complexity of $O(nF^2)$. SGC [30] repeatedly perform multiplication of normalized adjacency matrix $\tilde{\mathbf{A}}$ and feature matrix \mathbf{X} in the precomputation phase, which requires $O(LmF)$ time. PPRGo [4] calculates approximate the Personalized PageRank (PPR) matrix $\sum_{\ell=0}^{\infty} \alpha(1-\alpha)^\ell \tilde{\mathbf{A}}^\ell$ by forward push algorithm [2] and then applies the PPR matrix to the feature matrix \mathbf{X} to derive the propagation matrix. Let ε denote the error threshold of the forward push algorithm, the precomputation cost is bounded by $O(\frac{m}{\varepsilon})$. A major drawback of PPRGo is that it takes $O(\frac{n}{\varepsilon})$ space to store the PPR matrix, rendering it infeasible on billion-scale graphs.

Table 1: Summary of time complexity for GNN training and inference.

Method	Precomputation	Training	Inference
GCN	-	$O(LmF + LnF^2)$	$O(LmF + LnF^2)$
GraphSAGE	-	$O(ns_n^L F + ns_n^{L-1} F^2)$	$O(ns_n^L F + ns_n^{L-1} F^2)$
FastGCN	-	$O(Lns_l F + LnF^2)$	$O(Lns_l F + LnF^2)$
LADIES	-	$O(Lns_l F + LnF^2)$	$O(Lns_l F + LnF^2)$
SGC	$O(LmF)$	$O(nF^2)$	$O(nF^2)$
PPRGo	$O(\frac{m}{\varepsilon})$	$O(nKF + LnF^2)$	$O(nKF + LnF^2)$
Cluster-GCN	$O(m)$	$O(LmF + LnF^2)$	$O(LmF + LnF^2)$
GraphSAINT	-	$O(LbdF + LnF^2)$	$O(LmF + LnF^2)$
GBP (This paper)	$O(LnF + L\frac{\sqrt{m \lg n}}{\varepsilon} F)$	$O(LnF^2)$	$O(LnF^2)$

Other related work. Another line of research devotes to attention model [29, 27, 22], where the adjacency matrix of each layer is replaced by a learnable attention matrix. GIN [31] studies the expressiveness of GNNs, and shows that GNNs are not better than the Weisfeiler-Lehman test in

distinguishing graph structures. GDC [17] proposes to replace the graph convolutional matrix $\tilde{\mathbf{A}}$ with a graph diffusion matrix, such as the Heat Kernel PageRank or the Personalized PageRank matrix. Mixhop [1] mixes higher-order information to learn a wider class of representations. JKNet [32] explores the relationship between node influence and random walk in GNNs. DropEdge [24] and PairNorm [38] focus on over-smoothing problem and improve performance on GCNs by increasing the number of layers. These works focus on the effectiveness of GNNs; Thus, they are orthogonal to this paper’s contributions.

3 Bidirectional Propagation Method

Generalized PageRank. To achieve high scalability, we borrow the idea of decoupling prediction and propagation from SGC [30] and APPNP [16]. In particular, we precompute the feature propagation with the following *Generalized PageRank* matrix [21]:

$$\mathbf{P} = \sum_{\ell=0}^L w_{\ell} \mathbf{T}^{(\ell)} = \sum_{\ell=0}^L w_{\ell} (\mathbf{D}^{r-1} \mathbf{A} \mathbf{D}^{-r})^{\ell} \cdot \mathbf{X}, \quad (2)$$

where $r \in [0, 1]$ is the convolution coefficient, w_{ℓ} ’s are the weights of different order convolution matrices that satisfy $\sum_{\ell=0}^{\infty} w_{\ell} \leq 1$, and $\mathbf{T}^{(\ell)} = (\mathbf{D}^{r-1} \mathbf{A} \mathbf{D}^{-r})^{\ell} \cdot \mathbf{X}$ denotes the ℓ -th step propagation matrix. After the feature propagation matrix \mathbf{P} is derived, we can apply a multi-layer neural network with mini-batch training to make the prediction. For example, for multi-class classification tasks, a two-layer GBP model makes prediction with $Y = \text{SoftMax}(\sigma(\mathbf{P} \mathbf{W}_1) \mathbf{W}_2)$ where \mathbf{W}_1 and \mathbf{W}_2 are the learnable weight matrices, and σ is the activation function.

We note that equation (2) can be easily generalized to various existing models. By setting $r = 0.5, 0$ and 1 , the convolution matrix $\mathbf{D}^{r-1} \mathbf{A} \mathbf{D}^{-r}$ represents the symmetric normalization adjacency matrix $\mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}$ [15, 30, 16], the transition probability matrix $\mathbf{A} \mathbf{D}^{-1}$ [11, 8, 37], and the reverse transition probability matrix $\mathbf{D}^{-1} \mathbf{A}$ [32], respectively. We can also manipulate the weights w_{ℓ} to simulate various diffusion processes as in [17]. However, we will mainly focus on two setups in this paper: 1) $w_{\ell} = \alpha(1 - \alpha)^{\ell}$ for some constant decay factor $\alpha \in (0, 1)$, in which case \mathbf{P} becomes the *Personalized PageRank* used in APPNP and PPRGo [16, 17, 4]; 2) $w_{\ell} = 0$ for $\ell = 0, \dots, L-1$ and $w_L = 1$, in which case \mathbf{P} degenerates to the L -th transition probability matrix in SGC [30].

Algorithm 1: Bidirectional Propagation Algorithm

Input: Graph G , level L , training set V_t , weight coefficients w_{ℓ} , convolutional coefficient r , threshold r_{max} , number of walks per node n_r , feature matrix $\mathbf{X}_{n \times F}$

Output: Embedding matrix $\hat{\mathbf{P}}_{n \times F}$

- 1 $\mathbf{S}^{(\ell)} \leftarrow \text{Sparse}(\mathbf{0}_{n \times n})$ for $\ell = 0, \dots, L$;
- 2 **for each node** $s \in V_t$ **do**
- 3 Generate n_r random walks from s , each of length L ;
- 4 **if** The j -th random walk visits node u at the ℓ -th step, $\ell = 0, \dots, L$, $j = 1, \dots, n_r$ **then**
- 5 $\mathbf{S}^{(\ell)}(s, u) += \frac{1}{n_r}$;
- 6 $\mathbf{Q}^{(\ell)}, \mathbf{R}^{(\ell)} \leftarrow \text{Sparse}(\mathbf{0}_{n \times F})$ for $\ell = 1, \dots, L$;
- 7 $\mathbf{Q}^{(0)} \leftarrow \mathbf{0}_{n \times F}$ and $\mathbf{R}^{(0)} \leftarrow \text{ColumnNormalized}(\mathbf{D}^{-r} \mathbf{X})$;
- 8 **for** ℓ from 0 to $L-1$ **do**
- 9 **for each** $u \in V$ and $k \in \{0, \dots, F-1\}$ with $|\mathbf{R}^{(\ell)}(u, k)| > r_{max}$ **do**
- 10 **for each** $v \in \mathcal{N}(u)$ **do**
- 11 $\mathbf{R}^{(\ell+1)}(v, k) += \frac{\mathbf{R}^{(\ell)}(u, k)}{d(v)}$;
- 12 $\mathbf{Q}^{(\ell)}(u, k) \leftarrow \mathbf{R}^{(\ell)}(u, k)$ and $\mathbf{R}^{(\ell)}(u, k) \leftarrow 0$;
- 13 $\mathbf{Q}^{(L)} \leftarrow \mathbf{R}^{(L)}$ and $\mathbf{R}^L \leftarrow \text{Sparse}(\mathbf{0}_{n \times F})$;
- 14 $\hat{\mathbf{P}} \leftarrow \sum_{\ell=0}^L w_{\ell} \cdot \mathbf{D}^r \cdot (\mathbf{Q}^{(\ell)} + \sum_{t=0}^{\ell} \mathbf{S}^{(\ell-t)} \mathbf{R}^{(t)})$;
- 15 **return** Embedding matrix $\hat{\mathbf{P}}_{n \times F}$;

3.1 The Bidirectional Propagation Algorithm

To reduce the time complexity, we propose approximating the Generalized PageRank matrix \mathbf{P} with a localized bidirectional propagation algorithm from both the training/testing nodes and the feature vectors. Similar techniques have been used for computing probabilities in Markov Process [3]. Algorithm 1 illustrates the pseudo-code of the Bidirectional Propagation algorithm. The algorithm proceeds in three phases: Monte-Carlo Propagation (lines 1-5), Reverse Push Propagation (lines 6-13), and the combining phase (line 14).

Monte Carlo Propagation from the training/testing nodes. We start with a simple **Monte-Carlo propagation** (Lines 1-5 in Algorithm 1) from the training/testing nodes to **estimate the transition probabilities**. Given a graph G and a training/testing node set V_t , we generate a number n_r of random walks from each node $s \in V_t$, and record $\mathbf{S}^{(\ell)}(s, u)$, the fraction of random walks that visit node u at the ℓ -th steps. Note that $\mathbf{S}^{(\ell)} \in \mathcal{R}_{n \times n}$ is a **sparse matrix with at most $|V_t| \cdot n_r$ non-zero entries**. Since each random walk is independently sampled, we have that $\mathbf{S}^{(\ell)}$ is an unbiased estimator for the ℓ -th transition probability matrix $(\mathbf{D}^{-1}\mathbf{A})^\ell$. We also note that $(\mathbf{D}^{r-1}\mathbf{A}\mathbf{D}^{-r})^\ell = \mathbf{D}^r (\mathbf{D}^{-1}\mathbf{A})^\ell \mathbf{D}^{-r}$, which means we can use $\mathbf{D}^r \mathbf{S}^{(\ell)} \mathbf{D}^{-r} \mathbf{X}$ as an unbiased estimator for the ℓ -th propagation matrix $\mathbf{T}^{(\ell)} = (\mathbf{D}^{r-1}\mathbf{A}\mathbf{D}^{-r})^\ell \mathbf{X}$. Consequently, $\sum_{\ell=0}^L w_\ell \mathbf{D}^r \mathbf{S}^{(\ell)} \mathbf{D}^{-r} \mathbf{X}$ serves as an unbiased estimator for the Generalized PageRank Matrix \mathbf{P} . However, this estimation requires a large number of random walks from each training node and thus is infeasible for fully-supervised training on large graphs. why

Reverse Push Propagation from the feature vectors. To **reduce the variance of the Monte-Carlo estimator**, we employ a deterministic Reverse Push Propagation (lines 6-13 in Algorithm 1) from the feature vectors. Given a feature matrix \mathbf{X} , the algorithm outputs two sparse matrices for each level $\ell = 0, \dots, L$: the reserve matrix $\mathbf{Q}^{(\ell)}$ that **represents the probability mass to stay at level ℓ** , and the residue matrix $\mathbf{R}^{(\ell)}$ that **represents the probability mass to be distributed beyond level ℓ** . We begin by setting the initial residue $\mathbf{R}^{(0)}$ as the degree normalized feature matrix $\mathbf{D}^{-r} \mathbf{X}$. We also perform column-normalization on $\mathbf{R}^{(0)}$ such that each dimension of $\mathbf{R}^{(0)}$ has the same L_1 -norm. Starting from level $\ell = 0$, if the absolute value of the residue entry $\mathbf{R}^{(\ell)}(u, k)$ exceeds a threshold r_{max} , we **increase the residue of each neighbor v at level $\ell + 1$ by $\frac{\mathbf{R}^{(\ell)}(u, k)}{d(v)}$ and transfer the residue of u to its reserve $\mathbf{Q}^{(\ell)}(u, k)$** . Note that by maintaining a list of residue entries $\mathbf{R}^{(\ell)}(u, k)$ that exceed the threshold r_{max} , the above push operation can be done without going through every entry in $\mathbf{R}^{(\ell)}$. Finally, we transfer the non-zero residue $\mathbf{R}^{(L)}(u, k)$ of each node u to its reserve at level L .

We will show that when Reverse Push Propagation terminates, the reserve matrix $\mathbf{Q}^{(\ell)}$ satisfies

$$\mathbf{T}^{(\ell)}(s, k) - d(s)^r \cdot (\ell + 1) \cdot r_{max} \leq (\mathbf{D}^r \mathbf{Q}^{(\ell)})(s, k) \leq \mathbf{T}^{(\ell)}(s, k) \quad (3)$$

for each training/testing node $s \in V_t$ and feature dimension k . Recall that $\mathbf{T}^{(\ell)} = (\mathbf{D}^{r-1}\mathbf{A}\mathbf{D}^{-r})^\ell \mathbf{X}$ is the ℓ -th propagation matrix. This property implies that we can use $\sum_{\ell=0}^L w_\ell \mathbf{D}^r \mathbf{Q}^{(\ell)}$ to approximate the Generalized PageRank matrix $\mathbf{P} = \sum_{\ell=0}^L w_\ell \mathbf{T}^{(\ell)}$. However, there are two **drawbacks** to this approximation. First, this **estimator is biased**, which could potentially hurt the performance of the prediction. Secondly, the Reverse Push Propagation does **not take advantage of the semi-supervised learning setting where the number of training nodes may be significantly less than the total number of nodes n** .

Combining Monte-Carlo and Reverse Push Propagation. Finally, we combine the results from the Monte-Carlo and Reverse Push Propagation to **form a more accurate unbiased estimator**. In particular, we use the following equation to derive an approximation of the ℓ -th propagation matrix $\mathbf{T}^{(\ell)} = (\mathbf{D}^{r-1}\mathbf{A}\mathbf{D}^{-r})^\ell \mathbf{X}$:

$$\hat{\mathbf{T}}^{(\ell)} = \mathbf{D}^r \cdot \left(\mathbf{Q}^{(\ell)} + \sum_{t=0}^{\ell} \mathbf{S}^{(\ell-t)} \mathbf{R}^{(t)} \right). \quad (4)$$

As we shall prove in Section 3.2, $\hat{\mathbf{T}}^{(\ell)}$ is an unbiased estimator for $\mathbf{T}^{(\ell)} = (\mathbf{D}^{r-1}\mathbf{A}\mathbf{D}^{-r})^\ell \mathbf{X}$. Consequently, we can use $\hat{\mathbf{P}} = \sum_{\ell=0}^L w_\ell \hat{\mathbf{T}}^{(\ell)} = \sum_{\ell=0}^L w_\ell \mathbf{D}^r \cdot \left(\mathbf{Q}^{(\ell)} + \sum_{t=0}^{\ell} \mathbf{S}^{(\ell-t)} \mathbf{R}^{(t)} \right)$ as an

unbiased estimator for the Generalized PageRank matrix \mathbf{P} defined in equation (2). To see why equation (4) is a better approximation than the naive estimator $\mathbf{D}^r \mathbf{S}^{(\ell)} \mathbf{D}^{-r} \mathbf{X}$, note that each entry in the residue matrix $\mathbf{R}^{(\ell)}$ is bounded by a small real number r_{max} , which means the variance of the Monte-Carlo estimator is reduced by a factor of r_{max} . It is also worth mentioning that the two matrices $\mathbf{S}^{(\ell-t)}$ and $\mathbf{R}^{(t)}$ are sparse, so the time complexity of the matrix multiplication only depends on their numbers of non-zero entries.

3.2 Analysis

We now analyze the time complexity and the approximation quality of the Bidirectional Propagation algorithm. Due to the space limit, we defer all proofs in this section to the appendix. Recall that $|V_t|$ is the number of training/testing nodes, n_r is the number of random walks per node, and r_{max} is the push threshold. We assume $\mathbf{D}^{-r} \mathbf{X}$ is column-normalized, as described in Algorithm 1. We first present a Lemma that bounds the time complexity of the Bidirectional Propagation algorithm.

Lemma 1. *The time complexity of Algorithm 1 is bounded by $O\left(L^2 |V_t| n_r F + \frac{LdF}{r_{max}}\right)$.*

Intuitively, the $L|V_t|n_r F$ term represents the time for the Monte-Carlo propagation, and the $\frac{LdF}{r_{max}}$ term is the time for the Reverse Push propagation. Next, we will show how to set the number of random walks n_r and the push threshold r_{max} to obtain a satisfying approximation quality. In particular, the following technical Lemma states that the Reverse Push Propagation maintains an invariant during the push process.

Lemma 2. *For any residue and reserve matrices $\mathbf{Q}^{(\ell)}, \mathbf{R}^{(\ell)}, \ell = 0, \dots, L$, we have*

$$\mathbf{T}^{(\ell)} = (\mathbf{D}^{r-1} \mathbf{A} \mathbf{D}^{-r})^\ell \mathbf{X} = \mathbf{D}^r \cdot \left(\mathbf{Q}^{(\ell)} + \sum_{t=0}^{\ell} (\mathbf{D}^{-1} \mathbf{A})^{\ell-t} \mathbf{R}^{(t)} \right). \quad (5)$$

We note that the only difference between equations (4) and (5) is that we replace $\mathbf{D}^{-1} \mathbf{A}$ with $\mathbf{S}^{(\ell)}$ in equation (4). Recall that $\mathbf{S}^{(\ell)}$ is an unbiased estimator for the ℓ -th transition probability matrix $(\mathbf{D}^{-1} \mathbf{A})^\ell$. Therefore, Lemma 2 ensures that $\hat{\mathbf{T}}^{(\ell)}$ is also an unbiased estimator for $\mathbf{T}^{(\ell)}$. Consequently, $\hat{\mathbf{P}} = \sum_{\ell=0}^L w_\ell \hat{\mathbf{T}}^{(\ell)}$ is an unbiased estimator for the propagation matrix \mathbf{P} . Finally, to minimize the overall time complexity of Algorithm 1 in Lemma 1, the general principle is to balance the costs of the Monte-Carlo and the Reverse Push propagations. In particular, we have the following Theorem.

Theorem 1. *By setting $n_r = O\left(\frac{1}{\varepsilon} \sqrt{\frac{d \log n}{|V_t|}}\right)$ and $r_{max} = O\left(\varepsilon \sqrt{\frac{d}{|V_t| \log n}}\right)$, Algorithm 1 produces an estimator $\hat{\mathbf{P}}$ of the propagation matrix \mathbf{P} , such that for any $s \in V_t$ and $k = 0, \dots, F-1$, the probability that $\left| \hat{\mathbf{P}}(s, k) - \mathbf{P}(s, k) \right| \leq d(s)^r \varepsilon$ is at least $1 - \frac{1}{n}$. The time complexity is bounded by $O\left(L \cdot \frac{\sqrt{L|V_t|d \log(nL)}}{\varepsilon} \cdot F\right)$.*

For fully-supervised learning, we have $|V_t| = n$ and thus the time complexity of GBP becomes $O\left(L \cdot \frac{\sqrt{Lm \log(nL)}}{\varepsilon} \cdot F\right)$. In practice, we can also make a trade-off between efficiency and accuracy by manipulating the push threshold r_{max} and the number of walks n_r .

Parallelism of GBP. The Bidirectional Propagation algorithm is embarrassingly parallelizable: We can generate the random walks on multiple nodes and perform Reverse Push on multiple features in parallel. After we obtain the Generalized PageRank matrix \mathbf{P} , it is trivially to construct mini-batches for training the neural networks.

4 Experiments

Datasets. We use seven open graph datasets with different size: three citation networks Cora, Citeseer and Pubmed [25], a Protein-Protein interaction network PPI [11], a customer interaction

Table 2: Dataset statistics.

Dataset	Task	Nodes	Edges	Features	Classes	Label rate
Cora	multi-class	2,708	5,429	1,433	7	0.052
Citeseer	multi-class	3,327	4,732	3,703	6	0.036
Pubmed	multi-class	19,717	44,338	500	3	0.003
PPI	multi-label	56,944	818,716	50	121	0.79
Yelp	multi-label	716,847	6,977,410	300	100	0.75
Amazon	multi-class	2,449,029	61,859,140	100	47	0.70
Friendster	multi-class	65,608,366	1,806,067,135	100 (random)	500	0.001

network Yelp [37], a co-purchasing networks Amazon [8] and a large social network Friendster [34]. Table 2 summarizes the statistics of the datasets. We first evaluate GBP’s performance for transductive semi-supervised learning on the three popular citation networks (Cora, Citeseer, and Pubmed). Then we compare GBP with scalable GNN methods three medium to large graphs PPI, Yelp, Amazon in terms of inductive learning ability. Finally, we present the first empirical study of transductive semi-supervised on billion-scale network Friendster.

Baselines and detailed setup. We adopt three state-of-the-art GNN methods GCN [15], GAT [29], GDC [17] and APPNP [16] as the baselines for evaluation on small graphs. We also use one state-of-the-art scalable GNN from each of the three categories: LADIES (layer sampling) [40], GraphSAINT (graph sampling) [37], SGC and PPRGo (linear model) [30, 4].

We implement GBP in PyTorch and C++, and employ initial residual connection [12] across the hidden layers to facilitate training. For simplicity, we use the Personalized PageRank weights ($w_\ell = \alpha(1 - \alpha)^\ell$ for some hyperparameter $\alpha \in (0, 1)$). As we shall see, this weight sequence generally achieves satisfying results on graphs with real-world features. On the Friendster dataset, where the features are random noises, we use both Personalized PageRank and transition probability ($w_L = 1, w_0 = \dots = w_{L-1} = 0$) for GBP. We set $L = 4$ across all datasets. Table 3 summaries other hyper-parameters of GBP on different datasets. We use Adam optimizer to train our model, with a maximum of 1000 epochs and a learning rate of 0.01. For a fair comparison, we use the officially released code of each baseline (see the supplementary materials for URL and commit numbers) and perform a grid search to tune hyperparameters for models. All the experiments are conducted on a machine with an NVIDIA TITAN V GPU (12GB memory), Intel Xeon CPU (2.20 GHz), and 256GB of RAM.

Table 3: Hyper-parameters of GBP. r_{max} is the Reverse Push Threshold, w is the number of random walks from the training nodes, w_ℓ is the weight sequence, r is the Laplacian parameter in the convolutional matrix $\mathbf{D}^{r-1}\mathbf{A}\mathbf{D}^{-r}$.

Data	Dropout	Hidden dimension	L_2	Batch size	r_{max}	w	w_ℓ	r
Cora	0.5	64	5e-4	16	1e-4	0	$\alpha(1 - \alpha)^\ell, \alpha = 0.1$	0.5
Citeseer	0.5	64	2e-1	64	1e-5	0	$\alpha(1 - \alpha)^\ell, \alpha = 0.15$	0.4
Pubmed	0.5	128	5e-4	16	1e-5	0	$\alpha(1 - \alpha)^\ell, \alpha = 0.05$	0.5
PPI	0.1	2048	-	2048	5e-7	0	$\alpha(1 - \alpha)^\ell, \alpha = 0.3$	0
Yelp	0.1	2048	-	30000	5e-7	0	$\alpha(1 - \alpha)^\ell, \alpha = 0.9$	0.3
Amazon	0.1	1024	-	100000	1e-7	0	$\alpha(1 - \alpha)^\ell, \alpha = 0.2$	0.2
Friendster (PPR)	0.1	128	-	2048	4e-8	10000	$\alpha(1 - \alpha)^\ell, \alpha = 0.1$	0.5
Friendster	0.1	128	-	2048	4e-8	10000	$\{0, 0, 0, 1\}$	0.5

Transductive learning on small graphs. Table 4 shows the results for the semi-supervised transductive node classification task on the three small standard graphs Cora, Citeseer, and Pubmed. Following [15], we apply the standard fixed training/validation/testing split with 20 nodes per class for training, 500 nodes for validation and 1,000 nodes for testing. For each method, we set the number of hidden layers to 2 and take the mean accuracy with the standard deviation after ten runs. We observe that GBP outperforms APPNP (and consequently all other baselines) across all datasets. For the scalable GNNs, SGC is outperformed by the vanilla GCN due to the simplification [30]. On the other hand, the results of LADIES and GraphSAINT are also not at par with the non-scalable GNNs

Table 4: Results on Cora, Citeseer and Pubmed.

Method	Cora	Citeseer	Pubmed
GCN	81.5 \pm 0.6	71.3 \pm 0.4	79.1 \pm 0.4
GAT	83.3 \pm 0.8	71.9 \pm 0.7	78.0 \pm 0.4
GDC	83.3 \pm 0.2	72.2 \pm 0.3	78.6 \pm 0.4
APPNP	83.3 \pm 0.3	71.4 \pm 0.6	80.1 \pm 0.2
SGC	81.0 \pm 0.1	71.8 \pm 0.1	79.0 \pm 0.1
LADIES	79.6 \pm 0.5	68.6 \pm 0.3	77.0 \pm 0.5
PPRGo	82.4 \pm 0.2	71.3 \pm 0.3	80.0 \pm 0.4
GraphSAINT	81.3 \pm 0.4	70.5 \pm 0.4	78.2 \pm 0.8
GBP	83.9 \pm 0.7	72.9 \pm 0.5	80.6 \pm 0.4

such as GAT or APPNP, which suggests that the sampling technique alone might not be sufficient to achieve satisfying performance on small graphs.

Table 5: Results of inductive learning with scalable GNNs.

		4-layer		6-layer		8-layer	
		F1-score	Time (s)	F1-score	Time (s)	F1-score	Time (s)
PPI	SGC	65.7 \pm 0.01	76	62.4 \pm 0.01	173	57.8 \pm 0.01	295
	LADIES	57.9 \pm 0.30	187	59.4 \pm 0.25	206	57.4 \pm 0.24	315
	PPRGo	61.5 \pm 0.13	866	61.1 \pm 0.02	1976	55.1 \pm 0.19	1080
	GraphSAINT	99.2 \pm 0.05	1291	99.4 \pm 0.03	1961	99.3 \pm 0.01	2615
	GBP	99.3 \pm 0.02	117	99.3 \pm 0.03	167	99.3 \pm 0.01	220
Yelp	SGC	41.5 \pm 0.21	43	36.8 \pm 0.33	70	34.8 \pm 0.52	92
	LADIES	27.3 \pm 0.56	34	28.5 \pm 0.97	39	30.0 \pm 0.32	51
	PPRGo	64.0 \pm 0.16	550	63.7 \pm 0.71	1215	63.4 \pm 0.49	1665
	GraphSAINT	64.7 \pm 0.08	712	62.0 \pm 0.10	996	59.1 \pm 0.35	1298
	GBP	65.4 \pm 0.03	19	65.5 \pm 0.03	30	65.4 \pm 0.05	37
Amazon	SGC	90.4 \pm 0.01	233	89.9 \pm 0.03	284	89.7 \pm 0.03	342
	LADIES	85.4 \pm 0.14	734	85.2 \pm 0.20	784	84.6 \pm 0.09	1421
	PPRGo	83.3 \pm 0.51	2775	83.3 \pm 0.09	5206	81.6 \pm 0.22	9300
	GraphSAINT	91.5 \pm 0.01	957	91.3 \pm 0.05	1228	91.4 \pm 0.05	2618
	GBP	91.5 \pm 0.01	225	91.5 \pm 0.01	243	91.6 \pm 0.01	300

Inductive learning on medium to large graphs. Table 5 reports the F1-score and running time (precomputation + training) of each method with various depths on three large datasets PPI, Yelp, and Amazon. For each dataset, we set the hidden dimension to be the same across all methods: 2048(PPI), 2048(Yelp), and 1024(Amazon). We use “fixed-partition” splits for each dataset, following [37, 8] (see the supplementary materials for further details). The critical observation is that GBP can achieve comparable performance as GraphSAINT does, with 5-10x less running time. This demonstrates the superiority of GBP’s sub-linear time complexity. For PPRGo, it has a longer running time than other methods because of its expensive feature propagation per epoch. On the other hand, SGC and LADIES are also able to run faster than GraphSAINT; However, these two models’ accuracy is not comparable to that of GraphSAINT and GBP.

Figure 1 shows the convergence rate of each method, in which the time for data loading, pre-processing, validation set evaluation, and model saving are excluded. We observe that the convergence rate of GBP and SGC is much faster than that of LADIES and GraphSAINT, which is a benefit from decoupling the feature propagation and the neural networks.

Transductive semi-supervised learning on billion-scale graph Friendster. Finally, we perform the first empirical evaluation of scalable GNNs on a billion-scale graph Friendster. We extracted the top-500 ground-truth communities from [34] and use the community ids as the labels of each node. Note that one node may belong to multiple communities, in which case we pick the largest community as its label. The goal is to perform multi-class classification with only the graph structural information. This setting has been adapted in various works on community detection [19, 18, 35]. For each node, we generate a sparse random feature by randomly set one entry to be 1 in an d -dimensional all-zero vector. Note that even with a random feature matrix, GNNs are still able to extract structural

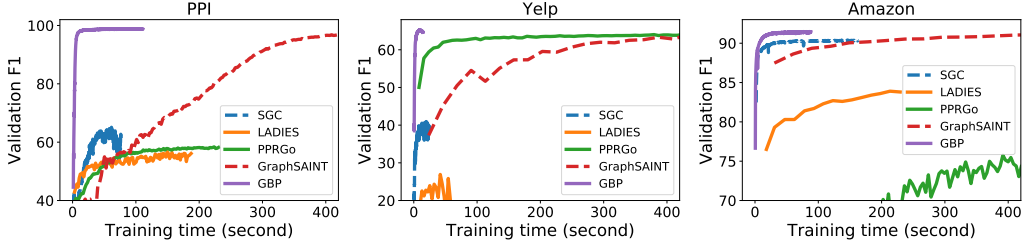


Figure 1: Convergence curves of 4-layer models.

information to perform the prediction [33]. Among the labeled nodes, we use 50,000 nodes (100 from each class) for training, 15,982 for validation, and 25,000 for testing. Table 6 report the running time and F1-score of each method with feature dimension $F = 10, 40, 70$ and 100. We omit GraphSAINT and LADIES as they run out of the 256 GB memory even with the dimension d set to 10. We first observe that both GBP and SGC can capture the structural information with random features, while PPRGo and GBP(PPR) fail to converge. This is because **Personalized PageRank emphasizes each node’s original feature** (with w_0 being the maximum weight among w_0, \dots, w_L) and, yet, the original features of Friendster are random noises. We also point out that PPRGo starts to converge and achieves an F1-score of 0.15 in 4500 seconds when the feature dimension is increased to 10000. On the other hand, we observe that GBP can achieve a significantly higher F1-score with less running time. Notably, on this 500-class classification task, GBP is able to achieve an F1-score of 0.79 with less than half an hour.

Table 6: Results for semi-supervised learning on Friendster.

Dimension	10		40		70		100	
F1-score / Time	F1-score	Time	F1-score	Time	F1-score	Time	F1-score	Time
SGC	2.0 ± 0.27	1130	12.9 ± 0.01	2930	27.1 ± 0.01	4549	40.2 ± 0.01	6379
PPRGo	1.6 ± 0.01	-	1.6 ± 0.01	-	1.6 ± 0.01	-	1.6 ± 0.01	-
GBP(PPR)	1.6 ± 0.01	-	1.6 ± 0.01	-	7.3 ± 0.20	-	7.3 ± 0.12	-
GBP	7.5 ± 0.10	757	26.6 ± 0.04	863	50.3 ± 0.44	1392	79.7 ± 0.32	1849

5 Conclusion

This paper presents GBP, a scalable GNN based on localized bidirectional propagation. Theoretically, GBP is the first method that achieves sub-linear time complexity for precomputation, training, and inference. The bidirectional propagation process computes a Generalized PageRank matrix that can express various existing graph convolutions. Extensive experiments on real-world graphs show that GBP obtains significant improvement over the state-of-the-art methods in terms of efficiency and performance. Furthermore, GBP is the first method that can scale to billion-edge networks on a single machine. For future work, an interesting direction is to extend GBP to heterogeneous networks.

Broader Impact

The proposed GBP algorithm addresses the challenge of scaling GNNs on large graphs. We consider this algorithm a general technical and theoretical contribution, without any foreseeable specific impacts. For applications in bioinformatics, computer vision, and natural language processing, applying the GBP algorithm may improve the scalability of existing GNN models. We leave the exploration of other potential impacts to future work.

Acknowledgments and Disclosure of Funding

Ji-Rong Wen was supported by National Natural Science Foundation of China (NSFC) No.61832017, and by Beijing Outstanding Young Scientist Program NO. BJJWZYJH012019100020098. Zhewei Wei was supported by NSFC No. 61972401 and No. 61932001, by the Fundamental Research Funds for the Central Universities and the Research Funds of Renmin University of China under Grant 18XNLG21, and by Alibaba Group through Alibaba Innovative Research Program. Ye Yuan was supported by NSFC No. 61932004 and No. 61622202, and by FRFCU No. N181605012. Xiaoyong Du was supported by NSFC No. U1711261.

References

- [1] ABU-EL-HAJJA, S., PEROZZI, B., KAPOOR, A., ALIPOURFARD, N., LERMAN, K., HARUTYUNYAN, H., STEEG, G. V., AND GALSTYAN, A. Mixhop: Higher-order graph convolutional architectures via sparsified neighborhood mixing. In *ICML* (2019).
- [2] ANDERSEN, R., CHUNG, F. R. K., AND LANG, K. J. Local graph partitioning using pagerank vectors. In *FOCS* (2006), IEEE Computer Society, pp. 475–486.
- [3] BANERJEE, S., AND LOFGREN, P. Fast bidirectional probability estimation in markov models. In *NIPS* (2015), pp. 1423–1431.
- [4] BOJCHEVSKI, A., KLICPERA, J., PEROZZI, B., KAPOOR, A., BLAIS, M., RÓZEMBERCZKI, B., LUKASIK, M., AND GÜNNEMANN, S. Scaling graph neural networks with approximate pagerank. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (New York, NY, USA, 2020), ACM.
- [5] CHEN, J., MA, T., AND XIAO, C. Fastgcn: Fast learning with graph convolutional networks via importance sampling. In *ICLR* (2018).
- [6] CHEN, J., ZHU, J., AND SONG, L. Stochastic training of graph convolutional networks with variance reduction. In *ICML* (2018).
- [7] CHEN, Z., WEI, X., WANG, P., AND GUO, Y. Multi-label image recognition with graph convolutional networks. In *CVPR* (2019).
- [8] CHIANG, W., LIU, X., SI, S., LI, Y., BENGIO, S., AND HSIEH, C. Cluster-gcn: An efficient algorithm for training deep and large graph convolutional networks. In *KDD* (2019), ACM, pp. 257–266.
- [9] CHUNG, F. R. K., AND LU, L. Survey: Concentration inequalities and martingale inequalities: A survey. *Internet Math.* 3, 1 (2006), 79–127.
- [10] FOUT, A., BYRD, J., SHARIAT, B., AND BEN-HUR, A. Protein interface prediction using graph convolutional networks. In *NeurIPS* (2017), pp. 6530–6539.
- [11] HAMILTON, W. L., YING, R., AND LESKOVEC, J. Inductive representation learning on large graphs. In *NeurIPS* (2017).
- [12] HE, K., ZHANG, X., REN, S., AND SUN, J. Deep residual learning for image recognition. In *CVPR* (2016), pp. 770–778.
- [13] HOSSEINI, B., MONTAGNE, R., AND HAMMER, B. Deep-aligned convolutional neural network for skeleton-based action recognition and segmentation. In *Data Science and Engineering* (2020).
- [14] KARYPIS, G., AND KUMAR, V. A fast and high quality multilevel scheme for partitioning irregular graphs. *SIAM J. Scientific Computing* 20, 1 (1998), 359–392.
- [15] KIPF, T. N., AND WELLING, M. Semi-supervised classification with graph convolutional networks. In *ICLR* (2017).
- [16] KLICPERA, J., BOJCHEVSKI, A., AND GÜNNEMANN, S. Predict then propagate: Graph neural networks meet personalized pagerank. In *ICLR* (2019).
- [17] KLICPERA, J., WEISSENBERGER, S., AND GÜNNEMANN, S. Diffusion improves graph learning. In *NeurIPS* (2019), pp. 13333–13345.
- [18] KLOSTER, K., AND GLEICH, D. F. Heat kernel based community detection. In *KDD* (2014), ACM, pp. 1386–1395.
- [19] LESKOVEC, J., LANG, K. J., AND MAHONEY, M. W. Empirical comparison of algorithms for network community detection. In *WWW* (2010), ACM, pp. 631–640.
- [20] LI, C., AND GOLDWASSER, D. Encoding social information with graph convolutional networks for political perspective detection in news media. In *ACL* (2019).
- [21] LI, P., CHIEN, I. E., AND MILENKOVIC, O. Optimizing generalized pagerank methods for seed-expansion community detection. In *NeurIPS* (2019), H. M. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. B. Fox, and R. Garnett, Eds., pp. 11705–11716.
- [22] LU, H., HUANG, S. H., YE, T., AND GUO, X. Graph star net for generalized multi-task learning. *CoRR abs/1906.12330* (2019).

- [23] QIU, J., TANG, J., MA, H., DONG, Y., WANG, K., AND TANG, J. Deepinf: Social influence prediction with deep learning. In *KDD* (2018), ACM, pp. 2110–2119.
- [24] RONG, Y., HUANG, W., XU, T., AND HUANG, J. Dropedge: Towards deep graph convolutional networks on node classification. In *ICLR* (2020).
- [25] SEN, P., NAMATA, G., BILGIC, M., GETOOR, L., GALLAGHER, B., AND ELIASSI-RAD, T. Collective classification in network data. *AI Magazine* 29, 3 (2008), 93–106.
- [26] SHANG, J., XIAO, C., MA, T., LI, H., AND SUN, J. Gamenet: Graph augmented memory networks for recommending medication combination. In *AAAI* (2019).
- [27] THEKUMPARAMPIL, K. K., WANG, C., OH, S., AND LI, L. Attention-based graph neural network for semi-supervised learning. *CoRR abs/1803.03735* (2018).
- [28] TONG, P., ZHANG, Q., AND YAO, J. Leveraging domain context for question answering over knowledge graph. In *Data Science and Engineering* (2019).
- [29] VELIČKOVIĆ, P., CUCURULL, G., CASANOVA, A., ROMERO, A., LIÒ, P., AND BENGIO, Y. Graph Attention Networks. *ICLR* (2018).
- [30] WU, F., SOUZA, A., ZHANG, T., FIFTY, C., YU, T., AND WEINBERGER, K. Simplifying graph convolutional networks. In *ICML* (2019), pp. 6861–6871.
- [31] XU, K., HU, W., LESKOVEC, J., AND JEGELKA, S. How powerful are graph neural networks? In *ICLR* (2019).
- [32] XU, K., LI, C., TIAN, Y., SONOBE, T., KAWARABAYASHI, K., AND JEGELKA, S. Representation learning on graphs with jumping knowledge networks. In *ICML* (2018).
- [33] YANG, C., XIAO, Y., ZHANG, Y., SUN, Y., AND HAN, J. Heterogeneous network representation learning: Survey, benchmark, evaluation, and beyond. *arXiv preprint arXiv:2004.00216* (2020).
- [34] YANG, J., AND LESKOVEC, J. Defining and evaluating network communities based on ground-truth. In *ICDM* (2012), IEEE Computer Society, pp. 745–754.
- [35] YANG, R., XIAO, X., WEI, Z., BHOWMICK, S. S., ZHAO, J., AND LI, R. Efficient estimation of heat kernel pagerank for local clustering. In *SIGMOD Conference* (2019), ACM, pp. 1339–1356.
- [36] YING, R., HE, R., CHEN, K., EKSOMBATCHAI, P., HAMILTON, W. L., AND LESKOVEC, J. Graph convolutional neural networks for web-scale recommender systems. In *KDD* (2018), ACM, pp. 974–983.
- [37] ZENG, H., ZHOU, H., SRIVASTAVA, A., KANNAN, R., AND PRASANNA, V. GraphSAINT: Graph sampling based inductive learning method. In *ICLR* (2020).
- [38] ZHAO, L., AND AKOGLU, L. Pairnorm: Tackling oversmoothing in gnns. In *ICLR* (2020).
- [39] ZHAO, L., PENG, X., TIAN, Y., KAPADIA, M., AND METAXAS, D. N. Semantic graph convolutional networks for 3d human pose regression. In *CVPR* (2019), pp. 3425–3435.
- [40] ZOU, D., HU, Z., WANG, Y., JIANG, S., SUN, Y., AND GU, Q. Layer-dependent importance sampling for training deep and large graph convolutional networks. In *NeurIPS* (2019), pp. 11247–11256.

A Proofs

We need the following Chernoff Bound for bounded i.i.d. random variables.

Lemma 3 (Chernoff Bound [9]). *Consider a set $\{x_i\}$ ($i \in [1, n_r]$) of i.i.d. random variables with mean μ and $x_i \in [0, r]$, we have*

$$\Pr \left[\left| \frac{1}{n_r} \sum_{i=1}^{n_r} x_i - \mu \right| \geq \varepsilon \right] \leq \exp \left(-\frac{n_r \cdot \varepsilon^2}{r \left(\frac{2}{3}\varepsilon + 2\mu \right)} \right). \quad (6)$$

A.1 Proof of Lemma 1

We note that the Bidirectional Propagation process shown in Algorithm 1 can be divided into three phases: the Monte Carlo Propagation phase (Line 1-5 in Algorithm 1), the Reverse Push Propagation phase (Line 6-13) and the combination phase (Line 14). In the combination phase, there are at most $O(|V_t| \cdot n_r)$ non-zero entries in each derived $\mathbf{S}^{(\ell)}$. Thus, the time cost to compute $\sum_{\ell=0}^L w_\ell \sum_{t=0}^{\ell} \mathbf{S}^{(\ell-t)} \mathbf{R}^{(t)}$ can be bounded by $O(L^2 |V_t| n_r F)$. Additionally, adding $\sum_{\ell=0}^L w_\ell \sum_{t=0}^{\ell} \mathbf{S}^{(\ell-t)} \mathbf{R}^{(t)}$ to $\sum_{\ell=0}^L \mathbf{Q}^{(\ell)}$ costs $O(L |V_t| F)$. Hence, the time cost of the combination phase is bounded by $O(L^2 |V_t| n_r F)$.

In the Monte-Carlo Propagation phase of Algorithm 1, we first generate n_r random walks of length L for each training/testing node $s \in V_t$ to estimate the ℓ -th transition probability matrix $\mathbf{S}^{(\ell)}$, $\ell = 0, \dots, L$. Since the number of training/testing nodes is $|V_t|$, the total cost is bounded by $O(L |V_t| n_r)$.

Furthermore, the time cost of the Reverse Push Propagation can be bounded by $\frac{LdF}{r_{max}}$, where d denotes the average degree of the given graph. Before proving this cost, we first introduce Lemma 4 to bound the sum of each column in $\mathbf{Q}^{(\ell)}$.

Lemma 4. *After the termination of Algorithm 1, the sum of the k -th column ($\forall k \in \{0, \dots, F-1\}$) in $\mathbf{Q}^{(\ell)}$ ($\forall \ell \in \{0, \dots, L\}$) can be bounded as $\mathbb{E} [\sum_{u \in V} |\mathbf{Q}^{(\ell)}(u, k)|] \leq 1$.*

Proof. As we shall prove in Lemma 2, any reserve matrix $\mathbf{Q}^{(\ell)}$ ($\forall \ell = 0, \dots, L$) can be bounded as:

$$\mathbf{D}^r \cdot \mathbf{Q}^{(\ell)} \leq \mathbf{T}^{(\ell)} = (\mathbf{D}^{r-1} \mathbf{A} \mathbf{D}^{-r})^\ell \cdot \mathbf{X} = \mathbf{D}^r \cdot (\mathbf{D}^{-1} \mathbf{A})^\ell \cdot (\mathbf{D}^{-r} \mathbf{X}),$$

following $\mathbf{Q}^{(\ell)} \leq (\mathbf{D}^{-1} \mathbf{A})^\ell \cdot (\mathbf{D}^{-r} \mathbf{X})$. Hence, for $\forall k \in \{0, \dots, F-1\}$ and $\forall u \in V$, we have:

$$\mathbf{Q}^{(\ell)}(u, k) \leq \left[(\mathbf{D}^{-1} \mathbf{A})^\ell \cdot (\mathbf{D}^{-r} \mathbf{X}) \right](u, k) = \sum_{v \in V} (\mathbf{D}^{-1} \mathbf{A})^\ell(u, v) \cdot (\mathbf{D}^{-r} \mathbf{X})(v, k).$$

Consequently, the sum of each column in $\mathbf{Q}^{(\ell)}$ can be bounded by

$$\begin{aligned} \sum_{u \in V} |\mathbf{Q}^{(\ell)}(u, k)| &\leq \sum_{u \in V} \sum_{v \in V} |(\mathbf{D}^{-1} \mathbf{A})^\ell(u, v) \cdot (\mathbf{D}^{-r} \mathbf{X})(v, k)| \\ &= \sum_{v \in V} |(\mathbf{D}^{-r} \mathbf{X})(v, k)| \cdot \sum_{u \in V} (\mathbf{D}^{-1} \mathbf{A})^\ell(u, v) = \sum_{v \in V} |(\mathbf{D}^{-r} \mathbf{X})(v, k)| = 1. \end{aligned} \quad (7)$$

In the first equality, we apply the fact that each entry of $(\mathbf{D}^{-1} \mathbf{A})^\ell$ is nonzero. In the second equality, we apply the property that the sum of each column in the ℓ -hop reverse transition probability matrix $(\mathbf{D}^{-1} \mathbf{A})^\ell$ equals 1, where $\ell \in \{0, 1, \dots\}$. And in the last equality, we use the fact that $\mathbf{D}^{-r} \mathbf{X}$ is column normalized before the Reverse Push Propagation phase (line 7 in Algorithm 1). Thus, Lemma 4 follows. \square

Based on Lemma 4, we can further bound the time cost of the Reverse Propagation process. Recall that in the Reverse Propagation phase of Algorithm 1, we push the residue $\mathbf{R}^{(\ell)}(u, k)$ of node u to its

neighbors whenever $|\mathbf{R}^{(\ell)}(u, k)| > r_{max}$, $k = 0, \dots, F - 1$. By Lemma 4, for a given level ℓ and a given feature dimension k , the sum $\sum_{u \in V} |\mathbf{Q}^{(\ell)}(u, k)| \leq 1$. Because $\mathbf{Q}^\ell = \mathbf{R}^\ell$ before we empty \mathbf{R}^ℓ (Line 13 in Algorithm 1), we have $\sum_{u \in V} |\mathbf{R}^{(\ell)}(u, k)| \leq 1$. Thus, there are at most $1/r_{max}$ nodes with residues larger than r_{max} . For random features, the average cost for this push operation is d , the average degree of the graph. Consequently, the cost of Reverse Push for a given level ℓ and a given feature dimension k is $\frac{d}{r_{max}}$. Summing up $\ell = 0, \dots, L - 1$ and $k = 0, \dots, F - 1$, we can derive the total time cost of the Reverse Push Propagation phase as $\frac{LdF}{r_{max}}$.

By summing up the time cost of the three phases mentioned above, the time complexity of Algorithm 1 is bounded by $O\left(L^2|V_t|n_rF + \frac{LdF}{r_{max}}\right)$, and Lemma 1 follows.

A.2 Proof of Lemma 2

Let \mathcal{RHS} denote the right hand side of equation (5); We prove the Lemma by induction. Recall that in Algorithm 1, we initialize $\mathbf{Q}^{(t)} = 0$ and $\mathbf{R}^{(t)} = 0$ for $t = 0, \dots, \ell$, and $\mathbf{R}^{(0)} = \mathbf{D}^{-r}\mathbf{X}$. Consequently, we have

$$\mathcal{RHS} = \mathbf{D}^r (\mathbf{D}^{-1}\mathbf{A})^\ell \mathbf{R}^{(0)} = \mathbf{D}^r (\mathbf{D}^{-1}\mathbf{A})^\ell \mathbf{D}^{-r}\mathbf{X} = (\mathbf{D}^{r-1}\mathbf{A}\mathbf{D}^{-r})^\ell \mathbf{X} = \mathbf{T}^{(\ell)},$$

which is true by definition. Assuming Equation (5) holds at some stage, we will show that the invariant still holds after a push operation on node u . More specifically, let $\mathbf{I}_{uk} \in \mathcal{R}^{n \times F}$ denote the matrix with entry at (u, k) setting to 1 and the rest setting to zero. Consider a push operation on $u \in V$ and $k \in 0, \dots, F - 1$ with $|\mathbf{R}^{(t)}(u, k)| > r_{max}$. We have two cases:

(1) If $t \leq \ell - 1$, we have $\mathbf{R}^{(t)}$ is decremented by $\mathbf{R}^{(t)}(u, k) \cdot \mathbf{I}_{uk}$ and $\mathbf{R}^{(t+1)}$ is incremented by $\frac{\mathbf{R}^{(t)}(u, k)}{d(v)} \cdot \mathbf{I}_{vk}$ for each $v \in N(u)$. Consequently, we have

$$\begin{aligned} \mathcal{RHS} &= \mathbf{T}^{(\ell)} + \mathbf{D}^r \cdot (\mathbf{D}^{-1}\mathbf{A})^{\ell-t} (-\mathbf{R}^{(t)}(u, k) \cdot \mathbf{I}_{uk}) + \mathbf{D}^r (\mathbf{D}^{-1}\mathbf{A})^{\ell-t-1} \cdot \sum_{v \in \mathcal{N}(u)} \frac{\mathbf{R}^{(t)}(u, k)}{d(v)} \cdot \mathbf{I}_{vk} \\ &= \mathbf{T}^{(\ell)} + \mathbf{R}^{(t)}(u, k) \cdot \mathbf{D}^r (\mathbf{D}^{-1}\mathbf{A})^{\ell-t-1} \cdot \left(\sum_{v \in \mathcal{N}(u)} \frac{1}{d(v)} \cdot \mathbf{I}_{vk} - \mathbf{D}^{-1}\mathbf{A}\mathbf{I}_{uk} \right) \\ &= \mathbf{T}^{(\ell)} + \mathbf{R}^{(t)}(u, k) \cdot \mathbf{D}^r (\mathbf{D}^{-1}\mathbf{A})^{\ell-t-1} \mathbf{0} = \mathbf{T}^{(\ell)}. \end{aligned}$$

For the second last equation, we use the fact that $\sum_{v \in \mathcal{N}(u)} \frac{1}{d(v)} \cdot \mathbf{I}_{vk} = \mathbf{D}^{-1}\mathbf{A}\mathbf{I}_{uk}$.

(2) If $t = \ell$, we have $\mathbf{R}^{(\ell)}$ is decremented by $\mathbf{R}^{(\ell)}(u, k) \cdot \mathbf{I}_{uk}$ and $\mathbf{Q}^{(\ell)}$ is incremented by $\mathbf{R}^{(\ell)}(u, k) \cdot \mathbf{I}_{uk}$. Consequently, we have

$$\mathcal{RHS} = \mathbf{T}^{(\ell)} + \mathbf{D}^r \cdot (-\mathbf{R}^{(\ell)}(u, k) \cdot \mathbf{I}_{uk}) + \mathbf{D}^r \cdot (\mathbf{R}^{(\ell)}(u, k) \cdot \mathbf{I}_{uk}) = \mathbf{T}^{(\ell)} + \mathbf{D}^r \cdot \mathbf{0} = \mathbf{T}^{(\ell)}.$$

Therefore, the induction holds, and the Lemma follows.

A.3 Proof of Theorem 1

Recall that \mathbf{P} is the Generalized PageRank matrix that $\mathbf{P} = \sum_{\ell=0}^L w_\ell \mathbf{T}^{(\ell)}$. The weights w_ℓ satisfy $\sum_{\ell=0}^\infty w_\ell \leq 1$. This suggests that for any $s \in V_t$ and $k = 0, \dots, F - 1$, by showing

$$\left| \hat{\mathbf{T}}^{(\ell)}(s, k) - \mathbf{T}^{(\ell)}(s, k) \right| > d(s)^r \varepsilon \quad (8)$$

holds with probability at most $\frac{1}{nL}$, Algorithm 1 also achieves the desired accuracy. Note that here we also apply the union bound:

$$\begin{aligned} \Pr \left[\left| \hat{\mathbf{P}}(s, k) - \mathbf{P}(s, k) \right| \leq d(s)^r \varepsilon \right] &\geq \Pr \left[\bigcap_{\ell=0}^L \left| \hat{\mathbf{T}}^{(\ell)}(s, k) - \mathbf{T}^{(\ell)}(s, k) \right| \leq d(s)^r \varepsilon \right] \\ &= 1 - \Pr \left[\bigcup_{\ell=0}^L \left| \hat{\mathbf{T}}^{(\ell)}(s, k) - \mathbf{T}^{(\ell)}(s, k) \right| > d(s)^r \varepsilon \right] \geq 1 - \sum_{\ell=0}^L \Pr \left[\left| \hat{\mathbf{T}}^{(\ell)}(s, k) - \mathbf{T}^{(\ell)}(s, k) \right| > d(s)^r \varepsilon \right]. \end{aligned}$$

Recall that we define $\mathbf{T}^{(\ell)} = \mathbf{D}^r \cdot \left(\mathbf{Q}^{(\ell)} + \sum_{t=0}^{\ell} (\mathbf{D}^{-1} \mathbf{A})^{\ell-t} \mathbf{R}^{(t)} \right)$. And Lemma 2 ensures $\hat{\mathbf{T}}^{(\ell)} = \mathbf{D}^r \cdot \left(\mathbf{Q}^{(\ell)} + \sum_{t=0}^{\ell} \mathbf{S}^{(\ell-t)} \mathbf{R}^{(t)} \right)$ is an unbiased estimator for $\mathbf{T}^{(\ell)}$. Hence, we have:

$$\begin{aligned} \left| \hat{\mathbf{T}}^{(\ell)}(s, k) - \mathbf{T}^{(\ell)}(s, k) \right| &= d(s)^r \cdot \left| \sum_{t=0}^{\ell} \left(\mathbf{S}^{(\ell-t)} \mathbf{R}^{(t)} \right)(s, k) - \sum_{t=0}^{\ell} \left((\mathbf{D}^{-1} \mathbf{A})^{\ell-t} \cdot \mathbf{R}^{(t)} \right)(s, k) \right| \\ &= d(s)^r \cdot \left| \sum_{t=0}^{\ell} \sum_{u \in V} \mathbf{S}^{(\ell-t)}(s, u) \mathbf{R}^{(t)}(u, k) - \sum_{t=0}^{\ell} \sum_{u \in V} (\mathbf{D}^{-1} \mathbf{A})^{\ell-t}(s, u) \cdot \mathbf{R}^{(t)}(u, k) \right|. \end{aligned}$$

According to Algorithm 1, each entry in residue matrix $\mathbf{R}^{(\ell)}$ is bounded by r_{max} after the Reverse Push Propagation phase. It follows:

$$\left| \hat{\mathbf{T}}^{(\ell)}(s, k) - \mathbf{T}^{(\ell)}(s, k) \right| \leq d(s)^r \cdot \left| \sum_{t=0}^{\ell} \sum_{u \in V} \mathbf{S}^{(\ell-t)}(s, u) \cdot r_{max} - \sum_{t=0}^{\ell} \sum_{u \in V} (\mathbf{D}^{-1} \mathbf{A})^{\ell-t}(s, u) \cdot r_{max} \right|,$$

where $\mathbf{S}^{(\ell)}(s, u)$ records the fraction of n_r random walks starting from $s \in V_t$ visit node u at the ℓ -th steps. We can use the Chernoff Bound given in Lemma 3 to show $\left| \hat{\mathbf{T}}^{(\ell)}(s, k) - \mathbf{T}^{(\ell)}(s, k) \right|$ holds with high probability. To make use of the Chernoff Bound, we first need to bound r , the maximum value of the random variables $\sum_{t=0}^{\ell} \sum_{u \in V} \mathbf{S}^{(\ell-t)}(s, u) \cdot r_{max}$. We note that $\sum_{u \in V} \mathbf{S}^{(\ell-t)}(s, u) = \sum_{u \in V} \sum_{w=1}^{n_r} \frac{x_w^{(\ell-t)}(s, u)}{n_r} \leq 1$. Here $x_w^{(\ell-t)}(s, u)$ is an indicator variable. $x_w^{(\ell-t)}(s, u) = 1$ when the w -th random walks starting from s walk at node u at the $(\ell - t)$ -th step. Thus, we have $\sum_{t=0}^{\ell} \sum_{u \in V} \mathbf{S}^{(\ell-t)}(s, u) \cdot r_{max} \leq (\ell + 1) \cdot r_{max}$. Plugging $r = (\ell + 1) \cdot r_{max}$ into Lemma 3, we can derive:

$$\begin{aligned} \Pr \left[\left| \hat{\mathbf{T}}^{(\ell)}(s, k) - \mathbf{T}^{(\ell)}(s, k) \right| > d(s)^r \varepsilon \right] \\ \leq \Pr \left[\left| \sum_{t=0}^{\ell} \sum_{u \in V} \mathbf{S}^{(\ell-t)}(s, u) \cdot r_{max} - \sum_{t=0}^{\ell} \sum_{u \in V} (\mathbf{D}^{-1} \mathbf{A})^{\ell-t}(s, u) \cdot r_{max} \right| > \varepsilon \right] \leq \exp \left(\frac{-n_r \cdot \varepsilon^2}{(\ell + 1) r_{max} \left(\frac{2\varepsilon}{3} + 2\mu \right)} \right), \end{aligned}$$

where $\mu = \sum_{t=0}^{\ell} \sum_{u \in V} (\mathbf{D}^{-1} \mathbf{A})^{\ell-t}(s, u) \cdot r_{max}$. Applying the fact that for any $(\ell - t) \in \{0, 1, \dots\}$, $\sum_{u \in V} (\mathbf{D}^{-1} \mathbf{A})^{\ell-t}(s, u) = 1$, we have $\mu = \sum_{t=0}^{\ell} r_{max} \cdot \sum_{u \in V} (\mathbf{D}^{-1} \mathbf{A})^{\ell-t}(s, u) = (\ell + 1) \cdot r_{max}$. By setting $n_r = \frac{r_{max} \left(\frac{2\varepsilon L}{3} + 2r_{max} \cdot L^2 \right) \cdot \log(nL)}{\varepsilon^2}$, we can ensure:

$$\Pr \left[\left| \hat{\mathbf{T}}^{(\ell)}(s, k) - \mathbf{T}^{(\ell)}(s, k) \right| > d(s)^r \varepsilon \right] \leq \exp \left(\frac{-n_r \cdot \varepsilon^2}{(\ell + 1) r_{max} \left(\frac{2\varepsilon}{3} + 2\mu \right)} \right) \leq \frac{1}{nL},$$

In the experimental part, we set $L = 4$, which is much smaller than ε and r_{max} . Hence, $n_r = O \left(\frac{r_{max} \cdot \log(nL)}{\varepsilon^2} \right)$. According to Lemma 1, the time complexity of Algorithm 1 can be expressed as:

$$O \left(L \cdot \frac{L|V_t| r_{max} \cdot \log(nL)}{\varepsilon^2} \cdot F + L \cdot \frac{d}{r_{max}} \cdot F \right).$$

We observe that the above complexity is minimized when $L \cdot \frac{L|V_t| r_{max} \cdot \log(nL)}{\varepsilon^2} \cdot F = L \cdot \frac{d}{r_{max}} \cdot F$, which implies the optimal $r_{max} = \sqrt{\frac{\varepsilon^2 d}{L|V_t| \log(nL)}}$. Hence, the total time complexity of Algorithm 1 can be further bounded as $O \left(L \cdot \frac{\sqrt{L|V_t| d \log(nL)}}{\varepsilon} \cdot F \right)$, which follows Theorem 1.

A.4 Proof of inequality (3)

According to Lemma 2, the following Equation holds for any residue and reserve matrices $\mathbf{Q}^{(\ell)}, \mathbf{R}^{(\ell)}$, $\ell = 0, \dots, L$:

$$\mathbf{T}^{(\ell)} = \mathbf{D}^r \cdot \left(\mathbf{Q}^{(\ell)} + \sum_{t=0}^{\ell} (\mathbf{D}^{-1} \mathbf{A})^{\ell-t} \mathbf{R}^{(t)} \right). \quad (9)$$

Thus, for $\forall s \in V$ and $\forall k \in \{0, \dots, F-1\}$, we have

$$\mathbf{T}^\ell(s, k) \geq (\mathbf{D}^r \cdot \mathbf{Q}^\ell)(s, k). \quad (10)$$

On the other hand, when the Reverse Propagation process terminates, each residue entry satisfies $\mathbf{R}^{(\ell)}(s, k) \leq r_{max}$ for $\forall u \in V$ and $\forall k \in \{0, \dots, F-1\}$. Hence, we can derive:

$$\left(\sum_{t=0}^{\ell} (\mathbf{D}^{-1} \mathbf{A})^{\ell-t} \mathbf{R}^{(t)} \right)(s, k) = \sum_{t=0}^{\ell} \sum_{u \in V} (\mathbf{D}^{-1} \mathbf{A})^{\ell-t}(s, u) \cdot \mathbf{R}^{(t)}(u, k) \leq (\ell+1) \cdot r_{max}.$$

In the last inequality, we apply the fact that $\mathbf{R}^{(t)}(u, k) \leq r_{max}$ and $\sum_{u \in V} (\mathbf{D}^{-1} \mathbf{A})^{\ell-t}(s, u) = 1$. Plugging into Equation (9), we can further derive the upper bound of $\mathbf{T}^\ell(s, k)$ as below.

$$\mathbf{T}^{(\ell)}(s, k) \leq (\mathbf{D}^r \cdot \mathbf{Q}^\ell)(s, k) + d(s)^r \cdot (\ell+1) \cdot r_{max}. \quad (11)$$

By combining the inequality (10) and inequality (11), we have:

$$\mathbf{T}^{(\ell)}(s, k) - d(s)^r \cdot (\ell+1) \cdot r_{max} \leq (\mathbf{D}^r \cdot \mathbf{Q}^\ell)(s, k) \leq \mathbf{T}^{(\ell)}(s, k),$$

which follows inequality (3).

B Additional experimental results

B.1 Comparison of inference time

Figure 2 shows the inference time of each method. We observe that in terms of the inference time, the three linear models, SGC, PPRGo and GBP, have a significant advantage over the two sampling-based models, LADIES and GraphSAINT.

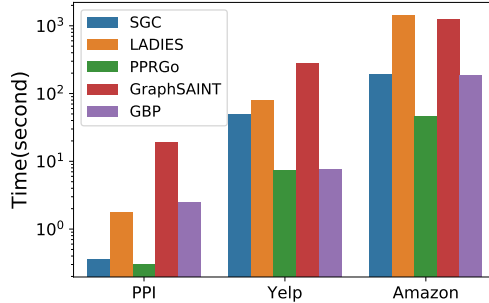


Figure 2: Inference time of 6-layers models on the entire test graph.

B.2 Additional details in experimental setup

Table 7 summarizes URLs and commit numbers of baseline codes.

Table 7: URLs of baseline codes.

Methods	URL	Commit
GCN	https://github.com/rustyls/pytorch_geometric	5692a8
GAT	https://github.com/rustyls/pytorch_geometric	5692a8
APPNP	https://github.com/rustyls/pytorch_geometric	5692a8
GDC	https://github.com/klicperajo/gdc	14333f
SGC	https://github.com/Tiiiger/SGC	6c450f
LADIES	https://github.com/acbull/LADIES	c7f987
PPRGo	https://github.com/TUM-DAML/pprgo_pytorch	d9f991
GraphSAINT	https://github.com/GraphSAINT/GraphSAINT	cd31c3