

Revisiting Graph Convolutional Network on Semi-Supervised Node Classification from an Optimization Perspective

Hongwei Zhang¹, Tijin Yan¹, Zejun Xie², Yuanqing Xia^{1*}, Yuan Zhang¹,

¹Beijing Institute of Technology, Beijing, China

²Renmin University of China, Beijing, China

Abstract

Graph convolutional networks (GCNs) have achieved promising performance on various graph-based tasks. However they suffer from over-smoothing when stacking more layers. In this paper, we present a quantitative study on this observation and develop novel insights towards the deeper GCN. First, we interpret the current graph convolutional operations from an optimization perspective and argue that over-smoothing is mainly caused by the naive first-order approximation of the solution to the optimization problem. Subsequently, we introduce two metrics to measure the over-smoothing on node-level tasks. Specifically, we calculate the fraction of the pairwise distance between connected and disconnected nodes to the overall distance respectively. Based on our theoretical and empirical analysis, we establish a universal theoretical framework of GCN from an optimization perspective and derive a novel convolutional kernel named GCN+ which has lower parameter amount while relieving the over-smoothing inherently. Extensive experiments on real-world datasets demonstrate the superior performance of GCN+ over state-of-the-art baseline methods on the node classification tasks.

1 Introduction

Graphs are ubiquitous in the real world, which can easily express various and complex relationships between objectives. In recent years, extensive studies have been conducted on deep learning methods for graph-structured data. There are several approaches on analyzing the graph, including network embedding (Perozzi, Al-Rfou, and Skiena 2014; Tang et al. 2015; Grover and Leskovec 2016), which only uses the graph structure, and graph neural networks (GNNs), which consider graph structure and node features simultaneously. GNNs have shown powerful ability on modeling the graph-structured data in a variety of graph learning tasks such as node classification (Gao, Wang, and Ji 2018; Hamilton, Ying, and Leskovec 2017; Yang, Cohen, and Salakhudinov 2016; Kipf and Welling 2016a; Veličković et al. 2018; Wu et al. 2019), link prediction (Zhang and Chen 2017, 2018; Cai and Ji 2020) and graph classification (Gilmer et al. 2017; Lee, Lee, and Kang 2019; Ma et al. 2019; Xu et al. 2018a; Ying et al. 2018b; Zhang et al. 2018). GNNs have also been applied to a range of applications, including social analysis (Qiu et al. 2018; Li and Goldwasser 2019), recommender systems (Ying et al. 2018a; Monti, Bronstein, and Bresson

2017), traffic prediction (Guo et al. 2019; Li et al. 2019b), drug discovery (Zitnik and Leskovec 2017) and fraud detection (Liu et al. 2019).

GNNs usually have different design paradigms, which include the spectral graph convolutional networks (Bruna et al. 2013; Defferrard, Bresson, and Vandergheynst 2016; Kipf and Welling 2016a), message passing framework (Gilmer et al. 2017; Hamilton, Ying, and Leskovec 2017), and neighbor aggregation via recurrent neural networks (Li et al. 2015; Dai et al. 2018). By using the idea of message passing framework, GNNs are to design various graph convolutional layers to update each node representation by aggregating the node representations from its neighbors.

However, most GNNs only consider the immediate neighbors for each node, which impedes their ability to extract the information of high-order neighbors. More layers usually lead to the performance degradation, which is caused by over-fitting and over-smoothing, of which the former is due to the increasing number of parameters when fitting a limited dataset whereas the latter is the inherent issue of the graph learning. How to make use of the high-order information of neighbors as well as achieving better performance remains a challenge. We need more insights to understand what GCN does and why over-smoothing occurs.

Several studies (Li, Han, and Wu 2018; Xu et al. 2018b; Klicpera, Bojchevski, and Günnemann 2019; Chen et al. 2020; Liu, Gao, and Ji 2020) have noticed over-smoothing, that is after multiple propagations, the final output of vanilla multi-layer GCN converges to a vector which only carries the information of the degree of graph and the node features are indistinguishable. Fig. 1 shows the node representations of vanilla multi-layer GCN on a small citation network *Cora*. We can observe that 2-layer GCN learns a meaningful embeddings which distinguish the different classes whereas more layers degrade the performance and lead to indistinguishable features.

Different from previous studies, we interpret the current graph convolutional operations from an optimization perspective, and argue that over-smoothing is mainly caused by the naive first-order approximation of the solution to the optimization problem. By solving it and applying the first-order approximation, we get the standard GCN kernel. This suggests that the original GCN kernel can be viewed as a simplified version of the solution. We argue that this sim-

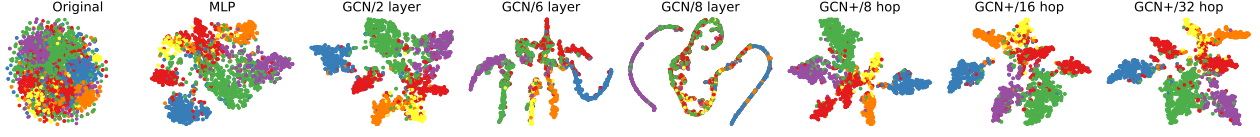


Figure 1: t-SNE Visualization of learned node representations, which include original features, MLP, different layers of GCN and different hops of GCN+ on *Cora*. Colors represent node classes.

plication loses necessary information which is crucial to tackle the over-smoothing to some extent. Based on this observation, two metrics are proposed to measure the smoothness of connected and disconnected pairwise node features respectively. Furthermore, we set three constraints: (a) the embedding learned by GCNs should not be too far off of the original features; (b) the connected nodes should have similar embeddings; (c) the disconnected nodes are assumed to have different embeddings.

As a result, we build a universal theoretical framework of GCN from an optimization perspective which smooths the node features and regularizes the (disconnected) node feature simultaneously. We consider two different cases of our framework, where the first case contains the current popular GCN (Kipf and Welling 2016a), SGC (Wu et al. 2019) and PPNP (Klicpera, Bojchevski, and Günnemann 2019), and the second case regularizes the pairwise distance of disconnected nodes.

The contributions of this work are summarized as follow:

- We provide a universal theoretical framework of GCN from an optimization perspective where the popular GCNs can be viewed as a special case of it. Furthermore, we derive a novel convolutional kernel named GCN+, which relieves the over-smoothing inherently and has lower parameter amount.
- We propose two quantitative metric to measure the smoothness and over-smoothness of the final nodes representations, which provides new insight to analyze the over-smoothing.
- We conduct extensive experiments on several public real-world datasets. Our results demonstrate the superior performance of GCN+ over state-of-the-art baseline methods.

2 Notations

Given an undirected graph $G = (V, E, X)$, V is node set with $|V| = n$, E is edge set. Let $A \in \mathbb{R}^{n \times n}$ denote the adjacency matrix, where $A_{ij} = 1$ if there is an edge between node i and node j otherwise 0. Let $D \in \mathbb{R}^{n \times n}$ denote the diagonal degree matrix where $D_{ii} = \sum_j A_{ij}$. Each node is associated with d features, and $X \in \mathbb{R}^{n \times d}$ is the feature matrix of nodes. each row of X is a signal defined over nodes. The graph Laplacian matrix is defined as $L = D - A$. Let $\tilde{A} = A + I$ and $\tilde{D} = D + I$ denote the adjacency and degree matrices of the self-loop graph respectively. We denote $\tilde{A}_{sym} = \tilde{D}^{-1/2} \tilde{A} \tilde{D}^{-1/2}$ and $\tilde{A}_{rw} = \tilde{D}^{-1} \tilde{A}$. Assume that each node v_i is associated with a class label $y_i \in Y$

where Y is a set of c classes. Let $N(v)$ denote the neighbors of v in graph, that is $N(v) = \{u \in V | \{u, v\} \in E\}$ and $\tilde{N}(v) = N(v) \cup \{v\}$. L' is the Laplacian matrix of the graph $G'(V', E', X)$, which is the complement of G , that means G' has the same nodes as G whereas if $\{u, v\} \in E$, then $\{u, v\} \notin E'$. Let A' and D' denote the corresponding adjacency and degree matrix respectively. We have $A' + A = J_n - I_n$ and $D' + D = (n - 1)I$ where J_n is a matrix whose element are all 1. Let $\text{num}(E)$ and $\text{num}(E')$ denote the numbers of edges in G and G' respectively, we have $\text{num}(E) + \text{num}(E') = \frac{n(n-1)}{2}$.

3 Perspectives of GCN

Here we provide three views to derive or understand the vanilla GCNs.

3.1 Spectral Graph Convolution

Bruna et al. (2013) define the spectral convolutions on graph by applying a filter g_θ in the Fourier domain to a graph signal. ChebNet (Defferrard, Bresson, and Vandergheynst 2016) suggests that the graph convolutional operation can be further approximated by the k -th order Chebyshev polynomial of Laplacian. Kipf and Welling (2016a) simplify the ChebNet and obtains a reduced version of ChebNet by the renormalization trick:

$$H^{(l+1)} = \sigma(\tilde{A}_{sym} H^{(l)} W^{(l)}) = \sigma(\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} H^{(l)} W^{(l)}), \quad (1)$$

where σ denote the activation function such as ReLU. $W^{(l)}$ is a layer-specific trainable weight matrix. $H^{(l)}$ is the feature matrix of l -th layer and $H^{(0)} = X$.

3.2 Message Passing

Message passing (Gilmer et al. 2017) means that a node on the graph aggregates the message from neighbors and update its embedding:

$$h_v^{(l)} = U_l \left(h_v^{(l-1)}, \sum_{u \in N(v)} M_l(h_u^{(l-1)}, h_v^{(l-1)}, e_{uv}) \right), \quad (2)$$

where $M_l(\cdot)$ and $U_l(\cdot)$ are message aggregation function and vertex update function, respectively. $h_v^{(l)}$ denotes the hidden state of node v at l -th layer, and e_{uv} is the edge features.

In this way, GCN layer can be decomposed into two steps, including the neighbors' message aggregation and update:

$$h_v^{(l)} = \sigma \left(W^{(l)} \sum_{u \in \tilde{N}(v)} \frac{h_u^{(l-1)}}{\sqrt{|N(v)||N(u)|}} \right). \quad (3)$$

Here a GCN layer can be viewed as a weighted average of all neighbors' message where the weighting is proportional to the inverse of the number of neighbors.

3.3 Graph Regularized Optimization

Let $\bar{X} \in \mathbb{R}^{n \times d}$ denote the final node embeddings matrix, and \bar{x}_i is the i -th row of \bar{X} . We consider the following optimization problem:

$$f = \min_{\bar{X}} \left(\sum_{i \in V} \|\bar{x}_i - x_i\|_D^2 + \alpha \sum_{\{i,j\} \in E} \|\bar{x}_i - \bar{x}_j\|_2^2 \right), \quad (4)$$

where $(x, y)_D = \sum_{i \in V} d(i)x(i)y(i)$, if $x = y$, we have $\|x\|_D = \sqrt{(x, x)_D}$.

The first term in the above optimization problem is the fitting constraint, which means the output features (also called embeddings) should not be too far off of the input features, while the second term is the smoothness constraint, which means the connected nodes should have similar embeddings. $\alpha > 0$ is a hyperparameter to balance the importance of two objections. It is worth noting that there is no limit to the specific transformation from X to \bar{X} .

Before solving the optimization problem, we have the following lemma.

Lemma 1 \tilde{A}_{rw} and \tilde{A}_{sym} always have the same eigenvalues $|\lambda| \leq 1$.

Corollary 1 $(I_n - \alpha\tilde{A}_{rw})$ and $(I_n - \alpha\tilde{A}_{sym})$ are invertible if $\alpha \in [0, 1)$.

Lemma 2 Given a graph with adjacency matrix A , the powers of A give the number of walks between any two vertices.

Corollary 2 A^k includes the information of high-order neighbors.

Next, we derive the closed-form solution of Eq. 4. Specifically, we rewrite Eq. 4 as

$$f = \min_{\bar{X}} \left(\text{Tr}((\bar{X} - X)(\bar{X} - X)^T \tilde{D}) + \alpha \text{Tr}(\bar{X}^T L \bar{X}) \right).$$

Differentiating f with respect to \bar{X} , we have

$$\frac{df}{d\bar{X}} = \tilde{D}(\bar{X} - X) + \alpha L \bar{X} = 0.$$

Notice Corollary 1, we have

$$\bar{X} = (1 - \mu)(I - \mu\tilde{A}_{rw})^{-1}X,$$

where $\mu = \frac{\alpha}{1+\alpha}$.

Actually, the solution is also the personalized PageRank (Page et al. 1999)'s limiting distribution. If we set $\mu = 0.5$, we get $\bar{X} = (2I - \tilde{A}_{rw})^{-1}X$, and $\tilde{A}_{rw}X$ is the first-order Taylor approximation. By replacing \tilde{A}_{rw} with \tilde{A}_{sym} , we get standard graph convolution kernel. In other words, we lose the information from high-order neighbors, which is contained in the error series of the Taylor expansion. (See Corollary 2).

In a nutshell, we obtain the well-known kernel or resemble form of the graph convolution from different ways.

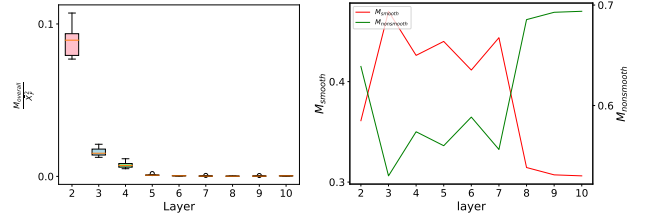


Figure 2: $M_{overall}$, M_{smooth} and $M_{non-smooth}$ of the output node embeddings of Vanilla GCN with increasing layers on *Cora*.

4 Over-smoothing in Vanilla Deep GCN

Neural network usually performs better when stack more layers while graph neural network does not benefit from the depth. On the contrary, more layers often result in significant degradation in performance.

Previous work illustrates the over-smoothing by computing the limiting distribution of A_k when $k \rightarrow \infty$. Actually, this is not identical with vanilla deep GCN, which contains non-linear transformation among different layers. Litter work considers the non-linearity in multi-layer GCN. Oono and Suzuki (2019) extend the linear analysis to the non-linearity firstly, which considers the ReLU activation function. They suggest that the node features of a k -layer GCNs will converge to a subspace and incur information loss, which makes the node feature indistinguishable.

At first, one main reason we introduce the deep architecture in GCN is that we want to use the long-range neighbor's information. We argue that vanilla deep GCN is not the correct way to capture this information. However, It does not mean that deep architecture is useless. Chen et al. (2020) and Liu, Gao, and Ji (2020) have shown that more layers can boost the performance of GCN on several datasets and tasks.

To quantify the over-smoothing in vanilla deep GCN, we compute the overall pairwise distance of node embeddings as follows:

$$\begin{aligned} M_{overall} &= \sum_{i,j \in V} \|\bar{x}_i - \bar{x}_j\|_2^2 = \sum_{\{i,j\} \in E} \|\bar{x}_i - \bar{x}_j\|_2^2 + \sum_{\{i,j\} \notin E} \|\bar{x}_i - \bar{x}_j\|_2^2 \\ &= \text{Tr}(\bar{X}^T L \bar{X}) + \text{Tr}(\bar{X}^T L' \bar{X}). \end{aligned}$$

Fig. 2(a) depicts the pairwise distance distribution of vanilla GCN with increasing layers on *Cora*. We can see that $M_{overall}$ decreases as the model goes deeper. Revisit the two parts of $M_{overall}$, we propose two fine quantitative metrics to measure the over-smoothing of graph representation.

$$\begin{aligned} M_{smooth} &= D_{smooth} / D_{overall}, \\ M_{non-smooth} &= D_{non-smooth} / D_{overall}, \end{aligned} \quad (5)$$

where

$$\begin{aligned} D_{smooth} &= \text{Tr}(\bar{X}^T L \bar{X}) / \text{num}(E), \\ D_{non-smooth} &= \text{Tr}(\bar{X}^T L' \bar{X}) / \text{num}(E'), \\ D_{overall} &= D_{smooth} + D_{non-smooth}. \end{aligned} \quad (6)$$

Here, $\text{num}(E)$ and $\text{num}(E')$ in the denominator are used to eliminate the impact of unbalanced edge numbers in G and G' . M_{smooth} measures the smoothness of the graph representation of connected pair nodes while $M_{non-smooth}$ measures the smoothness of the graph representation of disconnected pair nodes.

Fig. 2(b) compares M_{smooth} and $M_{non-smooth}$. We see that M_{smooth} contributes to quite a few parts of the overall distance, which seems counter-intuitive. We will discuss this two metrics of GCN+ in Section 6.3 again.

5 A General Framework of GCN

Recall the graph regularized optimization problem, we add a negative term to constrain the sum of distances between disconnected pairs as follow:

$$f = \min_{\bar{X}} \left(\sum_{i \in V} \|\bar{x}_i - x_i\|_D^2 + \alpha \sum_{\{i,j\} \in E} \|\bar{x}_i - \bar{x}_j\|_2^2 - \beta \sum_{\{i,j\} \notin E} \|\bar{x}_i - \bar{x}_j\|_2^2 \right),$$

where α and β are hyperparameters to balance the importance of the corresponding terms.

We consider two cases: $\beta = 0$ and $\beta \neq 0$.

5.1 Case 1: $\beta = 0$

In this situation, $\bar{X} = (1 - \mu)(I_n - \mu \tilde{A}_{rw})^{-1} X$ where $\mu = \frac{\alpha}{1+\alpha} \in (0, 1)$. Directly calculating such an intractable expression is not only computationally inefficient but also results in a dense $\mathbb{R}^{n \times n}$ matrix. It would lead to a high computational complexity and memory requirement when we apply such operator on large graphs. We can achieve linear computational complexity via power iteration.

We use \tilde{A} to denote \tilde{A}_{sym} and \tilde{A}_{rw} . Here we consider a more general expression $(1 - \mu)(I_n - \mu \tilde{A})^{-1} H$ where $H = f_\theta(X)$.

Theorem 1 $(I_n - \mu \tilde{A})$ is invertible. Consider the following iterative scheme

$$\begin{aligned} Z^{(0)} &= H, \\ Z^{(k)} &= \mu \tilde{A} Z^{(k-1)} + (1 - \mu) H, \end{aligned} \quad (7)$$

where $\mu \in (0, 1)$. When $k \rightarrow \infty$,

$$Z^{(\infty)} = (1 - \mu)(I_n - \mu \tilde{A})^{-1} H. \quad (8)$$

Proof Using corollary 1, we can see that $(I_n - \mu \tilde{A})$ is invertible. Combining the two equation of 7, we have

$$Z^{(k)} = \left(\mu^k \tilde{A}^k + (1 - \mu) \sum_{i=0}^{k-1} \mu^i \tilde{A}^i \right) H.$$

Notice that

$$\begin{aligned} \lim_{k \rightarrow \infty} \mu^k \tilde{A}^k &= 0, \\ \lim_{k \rightarrow \infty} \sum_{i=0}^{k-1} \mu^i \tilde{A}^i &= (I - \mu \tilde{A})^{-1}. \end{aligned} \quad (9)$$

Hence, the proof is finished.

Actually, the prevalent GCN, SGC and APPNP can be viewed as the special variant of Case 1.

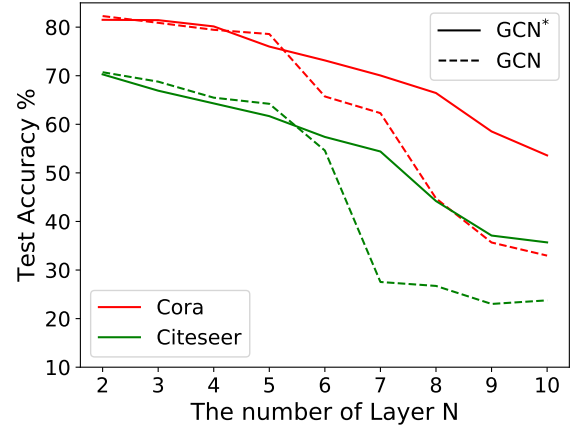


Figure 3: Performance comparison of vanilla deep GCN vs. GCN* with increasing layers on two small datasets.

5.2 Case 2: $\beta \neq 0$

In this situation, $\bar{X} = Q^{-1}$ if $Q = ((1 + \alpha + \beta)I - (\alpha + \beta)\tilde{D}^{-1}\tilde{A} - \beta n\tilde{D}^{-1} + \beta\tilde{D}^{-1}\mathbf{J}_n)$ is invertible when we choose a suitable β . We will introduce the conditions later.

First we use the first-order Taylor approximation of above convolutional kernel (GCN*) directly without any tricks such as Batch Normalization (Ioffe and Szegedy 2015) or residual connection (He et al. 2016) on two small citation datasets *Cora* and *Citeseer*. We compare the performance of the vanilla deep GCN and GCN* as the model layer increases. Fig. 3 shows the result of GCN and GCN*. Dashed lines illustrate the performance of GCN, which shows that deep GCN suffers from performance drop. We can see that the performance decay with GCN* kernel is much slower.

Oono and Suzuki (2019) have proved that the node feature of vanilla k -layer GCN will converges to an invariant subspace which only carry the information of the connected component and node degree. The convergence speed is proportional to the λ^k , where λ is the supremum of eigenvalue of \tilde{A} . In GCN*, $\lambda > 1$ (see the proof of Theorem 2), which implies that λ^k is large, thus the information loss and over-smoothing are relieved.

Although the modified graph kernel relieves over-smoothing to some extent, more layers do not boost the performance, which is not our focus. However the above result demonstrates that it is an efficient way to tackle the over-smoothing issue. We can achieves linear computational complexity via power iteration similar to Case 1.

Theorem 2 $(I_n - \mu \hat{A})$ is invertible when $\beta < \frac{1}{n}$ where $\mu = \frac{\alpha + \beta}{1 + \alpha + \beta}$ and $\hat{A} = \tilde{A} + \frac{\beta n \tilde{D}^{-1} - \beta \tilde{D}^{-1} \mathbf{J}_n}{\alpha + \beta}$. Consider the following iterative scheme

$$\begin{aligned} Z^{(0)} &= H, \\ Z^{(k)} &= \mu \hat{A} Z^{(k-1)} + (1 - \mu) H, \end{aligned} \quad (10)$$

where $\mu \in (0, 1)$. When $k \rightarrow \infty$,

$$Z^{(\infty)} = (1 - \mu)(I_n - \mu \hat{A})^{-1} H. \quad (11)$$

Proof Let $M = \frac{\beta n \bar{D}^{-1} - \beta \bar{D}^{-1} \mathbf{J}_n}{\alpha + \beta} = \frac{\beta \bar{D}^{-1}}{\alpha + \beta} (nI - \mathbf{J}_n)$. Note that $(nI - \mathbf{J}_n)$ has the largest eigenvalue n . Suppose that λ is the eigenvalue of M , we have $\lambda \leq \frac{\beta n}{\alpha + \beta}$. Then eigenvalue of \hat{A} is less than $1 + \frac{\beta n}{\alpha + \beta}$. $(I_n - \mu \hat{A})$ is invertible iff $\frac{1}{\mu}$ is not an eigenvalue of \hat{A} . Note that $\frac{1}{\mu} = \frac{1 + \alpha + \beta}{\alpha + \beta} = 1 + \frac{1}{\alpha + \beta}$, when $\beta < \frac{1}{n}$ we have $\frac{1}{\mu} > 1 + \frac{\beta n}{\alpha + \beta}$, hence $\frac{1}{\mu}$ cannot be an eigenvalue of \hat{A} and $(I_n - \mu \hat{A})$ is invertible. The proof of the iterative scheme follows the similar procedure of case 1 with a slight difference, as it is trivial, we omit the proof.

5.3 Why GCN+ relieve the over-smoothing?

We have no assumptions on the specific transformation from X to \bar{X} . In our implementation, the mathematical expression of GCN+ is defined as

$$\begin{aligned} Z^{(0)} &= H = \sigma(XW_1), \\ Z^{(k)} &= \mu \hat{A} Z^{(k-1)} + (1 - \mu)H, \\ X_{out} &= \text{softmax}(Z^{(k)}W_2), \end{aligned} \quad (12)$$

where $W_1 \in \mathbb{R}^{d \times m}$ and $W_2 \in \mathbb{R}^{m \times c}$ are learnable weight matrices, k is the dimension of the hidden layers.

We interpret the anti-oversmoothing of GCN+ from two ways. First, note that in the power iterative scheme, a fraction of initial node features H is always preserved in each iteration, which can be viewed as a flexible version of residual connection. In addition, we can also understand GCN+ from the frequency of graph signal. In Section 3.3, we have shown that the original GCN is corresponding to the first-order Taylor approximation of the optimization solution, that means we lost the high frequency part of the signal which contains the high-order information. Actually we omit the error series when we approximate the GCN.

Recall the current representative methods: DAGNN and JKNet, which shows promising improvement than the original GCN. The core formulas of them are as follows:

DAGNN:

$$Z = \text{stack}(H, Z^{(1)}, \dots, Z^{(k)}), \quad Z^{(k)} = \hat{A}^{(k)} H,$$

JKNet:

$$Z = \text{Aggr}(Z^{(1)}, \dots, Z^{(k)}),$$

where Aggr includes *Concatenation*, *Max-pooling* and *LSTM-attention*.

Actually, DAGNN and JKNet both make use of the information which from the immediate and high-order neighbors while GCN+ also benefit from this. Moreover, we provide the theoretical and empirical evidence of GCN+.

5.4 Parameters Amount

It is worth noting that the power iterative schemes are parameter-free in two versions of GCN+, which is similar to APPNP (Klicpera, Bojchevski, and Günnemann 2019). In particular, GCN+ ($\beta = 0$) adopts the same scheme as APPNP, where we re-implement it and achieve more impressive results.

Dataset	Nodes	Edges	Classes	Features	Metric
Cora	2708	5429	7	1433	Accuracy
Citeseer	3327	4732	7	2703	Accuracy
Pubmed	19717	44338	3	500	Accuracy
ogb-arxiv	169343	1166243	40	128	Accuracy
ogb-proteins	132534	39561252	112	8	ROC-AUC

Table 1: Dataset statistics.

6 Experiments

In this section, we evaluate the performance of GCN+ on several benchmark datasets against various graph neural networks on semi-supersized node classification tasks.

6.1 Experimental Setup

Datasets We conduct extensive experiments on the node-level tasks on two kinds commonly used networks: Planetoid: *Cora*, *CiteSeer*, *Pubmed* (Sen et al. 2008) and recent *Open Graph Benchmark* (OGB) (Hu et al. 2020): *ogb-arxiv*, *ogb-proteins*. The statistics of datasets are summarized in Table 1. It is worth nothing that OGB includes enormous challenging and large-scale datasets than Planetoid. We refer readers to (Hu et al. 2020) for more details about OGB datasets.

Implementations We choose the optimizer and hyperparameters of GNN models as follows. We use the Adam optimizer (Kingma and Ba 2014) to train all the GNN models with a maximum of 1500 epochs. We set the number of hidden units to 64 on *Cora*, *Citeseer* and *Pubmed*, 256 on *ogb-arxiv* and *ogb-proteins*. For SGC, we vary number of layer in $\{1, 2, \dots, 10, 15, \dots, 60\}$ and for GCN and GAT in $\{2, 4, \dots, 10, 15, \dots, 30\}$. For α in APPNP, we search it from $\{0.1, 0.2, 0.3, 0.4, 0.5\}$. For DAGNN and JKNet, we search layers from $\{2, 3, \dots, 10\}$. For learning rate, we choose from $\{0.001, 0.005, 0.01\}$. For dropout rate, we choose from $\{0.1, 0.2, 0.3, 0.4, 0.5\}$. We perform a grid search to tune the hyperparameters for other models based on the accuracy on the validation set. We run each experiment 10 times and report the average.

In practice, we use Pytorch (Paszke et al. 2019) and Pytorch Geometric (Fey and Lenssen 2019) for an efficient GPU-based implementation of GCN+. All experiments in this study are conducted on NVIDIA TITAN RTX 24GB GPU.

6.2 Comparison with SOTA

The evaluate metric of various datasets are listed in Table 1. Actually it is commonly used to evaluate the model by the community.

Planetoid We use standard fixed and random training/validation/testing splits. Specifically, we use 20 labeled nodes per class as the training set, 500 nodes as the validation set, and 1000 nodes as the test set for all models. For fixed split, we follow the experimental setup in (Yang, Cohen, and Salakhudinov 2016). We compare Multiplayer Perception (MLP), GCN (Kipf and Welling 2016a), GAT

model	Cora		Citeseer		Pubmed	
	Fixed	Random	Fixed	Random	Fixed	Random
MLP	61.6 ± 0.6	59.8 ± 2.4	61.0 ± 1.0	58.8 ± 2.2	74.2 ± 0.7	70.1 ± 2.4
GCN(Kipf and Welling 2016a)	81.3 ± 0.8	79.1 ± 1.8	71.1 ± 0.7	68.2 ± 1.6	78.8 ± 0.6	77.1 ± 2.7
GAT(Veličković et al. 2018)	83.1 ± 0.4	80.8 ± 1.6	70.8 ± 0.5	68.9 ± 1.7	79.1 ± 0.4	77.8 ± 2.1
SGC(Wu et al. 2019)	81.1 ± 0.5	80.4 ± 0.3	71.9 ± 0.3	71.8 ± 0.3	78.9 ± 0.0	77.8 ± 0.6
JKNet(Xu et al. 2018b)	80.7 ± 0.9	79.2 ± 0.9	70.1 ± 0.6	68.3 ± 1.8	78.1 ± 0.6	77.9 ± 0.9
APPNP(Klicpera, Bojchevski, and Günnemann 2019)	83.3 ± 0.5	81.9 ± 1.4	71.8 ± 0.4	69.8 ± 1.7	80.1 ± 0.2	79.5 ± 2.2
DAGNN(Liu, Gao, and Ji 2020)	84.4 ± 0.5	83.7 ± 1.4	73.3 ± 0.6	71.2 ± 1.4	80.5 ± 0.5	80.1 ± 1.7
GCN+($\beta = 0$)	85.2 ± 0.5	83.3 ± 1.1	73.3 ± 0.5	72.3 ± 0.7	80.4 ± 0.6	80.1 ± 0.6
GCN+($\beta \neq 0$)	85.6 ± 0.4	83.6 ± 1.3	73.5 ± 0.4	72.5 ± 0.9	80.5 ± 0.6	80.3 ± 0.7

Table 2: Summary of classification accuracy(%) on Planetoid datasets of semi-supervised node classification.

Dataset	<i>ogb-arxiv</i>	<i>ogb-proteins</i>
GCN	71.74 ± 0.29	72.51 ± 0.35
GraphSAGE	71.49 ± 0.25	77.68 ± 0.20
GCN+($\beta = 0$)	71.85 ± 0.23	78.63 ± 0.28
GCN+($\beta \neq 0$)	71.95 ± 0.28	79.07 ± 0.34

Table 3: Summary of classification performance(%) on OGB datasets. For *ogb-arxiv*, it indicates accuracy and for *ogb-proteins*, it indicates ROC-AUC.

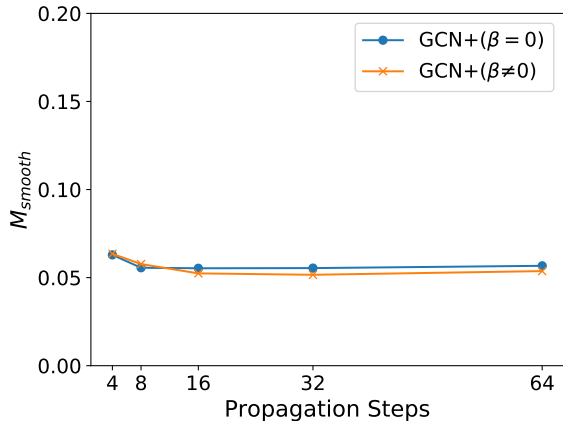


Figure 4: $M_{non-smooth}$ of GCN+ with increasing hops on Cora.

(Veličković et al. 2018), SGC (Wu et al. 2019), DAGNN (Liu, Gao, and Ji 2020) and APPNP (Klicpera, Bojchevski, and Günnemann 2019) with GCN+. Although DropEdge (Rong et al. 2019), PairNorm (Zhao and Akoglu 2019) are proposed to tackle over-smoothing issue recently, our baseline methods don't include them as they do not help to boost the performance on node classification task. Table 2 compares the average test accuracy of 10 runs for each model on Planetoid dataset. As shown, GCN+ outperforms better than the representative baselines. Note that the shallow model APPNP achieves better performance than GCN and GAT. Recent deeper model named DAGNN shows competitive result and robustness on these datasets and GCN+ per-

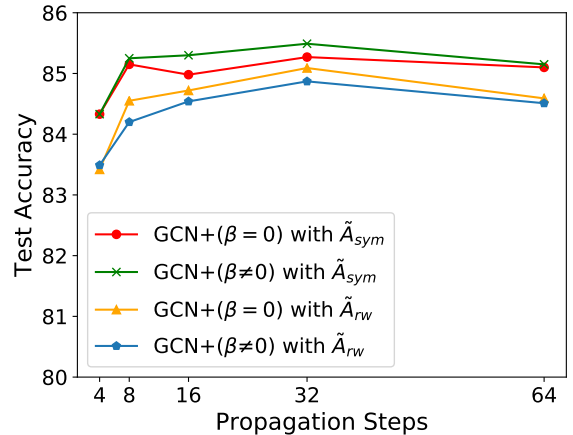


Figure 5: Performance comparison of different propagation matrices \tilde{A}_{sym} vs. \tilde{A}_{rw} in GCN+ with increasing hops on Cora.

forms slightly better than it.

OGB We adopt the setting of (Hu et al. 2020), which is more challenging and realistic. We consider the following representative models GCN (Kipf and Welling 2016a), GraphSAGE (Hamilton, Ying, and Leskovec 2017) and GCNII (Chen et al. 2020) as our baselines. In particular, we use the reported metric of the leaderboards of OGB team, which provide an open benchmark on several tasks and datasets.

Table 3 compares the average test accuracy/ROC-AUC on OGB datasets. As shown, GCN+ outperforms the GCN and GraphSAGE. It is clear that our proposed GCN+ outperform SOTA in two middle scale datasets.

In summary, GCN+ achieves superior performance on several benchmarks, which shows that considering the information of high-order neighbors makes sense and we need more reasonable way to deepen GCNs or make use of the high-order neighbors. Note that GCN+ ($\beta \neq 0$) is slightly better than GCN+ ($\beta = 0$) which is benefit from the third term of Eq. (7).

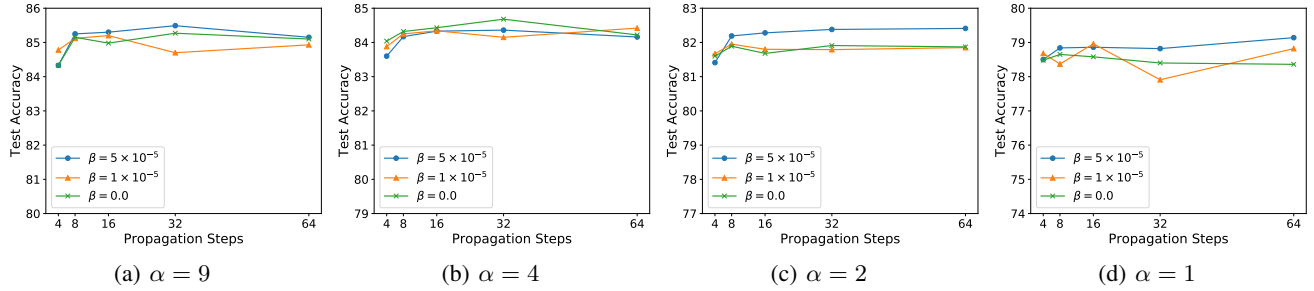


Figure 6: Test accuracy of different propagation steps and α on *Cora*.

6.3 Over-smoothing Analysis

We employ the two proposed metrics to measure the node embeddings learned by GCN+. The results on *Cora* are shown in Fig. 4. We can observe that as the number of hops increases, the M_{smooth} values nearly remains a small constant which is lower than vanilla deep GCN. This implies that GCN+ use the information of long-range neighbors and do not suffer from over-smoothing.

Fig. 1 also compares the final output embeddings of GCN+ with multiple hops, which shows different behaviors with GCN. GCN+ relieves the over-smoothing and learns the meaningful embeddings with the increasing hops.

6.4 Hyperparameter Analysis

In the previous sections, we use \tilde{A} to refer the \tilde{A}_{sym} and \tilde{A}_{rw} . Here we compare the different choices of propagation matrix \tilde{A} . Fig. 5 depicts the test accuracy achieved by varying the hops of different propagation matrices. The result illustrates that \tilde{A}_{sym} is slightly better than \tilde{A}_{rw} .

We consider three hyperparameter of GCN+, that is α , β and number of power iteration steps k . Fig. 6 compares the effect of these hyperparameters on *Cora*. We can see that $k = 16, 32$ is suitable and more steps does not boost the performance significantly. For *Cora*, when $\alpha = 9$ (that means the fraction of retained initial node features is 0.1.), GCN+ achieve the best performance. The value of α varies by different datasets. More results and details listed in the supplementary material.

7 Related Work

7.1 Graph Neural Networks

Graph neural networks (GNNs) have been extensively studied for the past years. There are different views on designing new architecture, including the spectral-based, spatial-based and other types, such as understand the GNN using dynamic system (Xhonneux, Qu, and Tang 2019). Numerous methods are proposed to model the graph-structure data and apply on a wide range of applications. Besides the GCNs, there are also other types of GNNs, such as attention-based GNN (Veličković et al. 2018) which use multiple attention to aggregate information from neighbors, autoencoder-based GNN (Kipf and Welling 2016b), which use a GCN encoder and decoder to learn meaningful embeddings, and dynamic

GNNs (Seo et al. 2018; Hajiramezanali et al. 2019; Yan et al. 2020) which learn the node embedding over time.

7.2 Deep GCN and Over-smoothing

Most GNNs are shallow models as deep architecture suffers from over-smoothing. Several studies explore deep GCNs. Xu et al. (2018b) introduce Jumping Knowledge Networks, which uses residual connection to combine the output of each layer. Klicpera, Bojchevski, and Günnemann (2019) use Personalized PageRank, which consider the information of root node to replace the graph convolution operator to solve the over-smoothing. DropEdge (Rong et al. 2019) suggests that randomly removing the edge of original graph impede over-smoothing. PairNorm (Zhao and Akoglu 2019) is another scheme which uses a normalization layer to scale the node features after the convolution layer. Li et al. (2019a) build on ideas from ResNet to train very deep GCNs. Li et al. (2020) further propose MsgNorm, which boosts the performance on several datasets. Yang et al. (2020) present NodeNorm to scale the node features. (Chen et al. 2020) propose a deep GCN models which use initial residual connection and identity mapping.

A few work analyzes the cause and behaviors of over-smoothing theoretically. Oono and Suzuki (2019) investigate the asymptotic behaviors of GCNs as the layer size tends to infinity and reveals the information loss in deep GCNs. Cai and Wang (2020) further extend analysis of (Oono and Suzuki 2019) from linear GNNs to the nonlinear architecture.

8 Conclusion

We summarize the existing different views on the mechanism of GCNs, which help us understand and design the graph convolutional kernel. We further provide a general optimization framework named GCN+. Based on this framework, we derive two forms of GCN+ and propose two metrics to measure the smoothness of output node representations. Extensive empirical studies on several real-world datasets demonstrate that GCN+ compares favorably to state of the art with a small amount of parameters. For future work, we will consider different optimization objectives which encode the graph structure and node features adaptively. As we do not limit the transformation from X to \tilde{X} , another reasonable formulas can be further explored.

References

- Bruna, J.; Zaremba, W.; Szlam, A.; and LeCun, Y. 2013. Spectral networks and locally connected networks on graphs. *ICLR*.
- Cai, C.; and Wang, Y. 2020. A Note on Over-Smoothing for Graph Neural Networks. *arXiv preprint arXiv:2006.13318*.
- Cai, L.; and Ji, S. 2020. A Multi-Scale Approach for Graph Link Prediction. In *AAAI*, 3308–3315.
- Chen, M.; Wei, Z.; Huang, Z.; Ding, B.; and Li, Y. 2020. Simple and Deep Graph Convolutional Networks. *arXiv preprint arXiv:2007.02133*.
- Dai, H.; Kozareva, Z.; Dai, B.; Smola, A.; and Song, L. 2018. Learning steady-states of iterative algorithms over graphs. In *International conference on machine learning*, 1106–1114.
- Defferrard, M.; Bresson, X.; and Vandergheynst, P. 2016. Convolutional neural networks on graphs with fast localized spectral filtering. In *Advances in neural information processing systems*, 3844–3852.
- Fey, M.; and Lenssen, J. E. 2019. Fast Graph Representation Learning with PyTorch Geometric. In *ICLR Workshop on Representation Learning on Graphs and Manifolds*.
- Gao, H.; Wang, Z.; and Ji, S. 2018. Large-Scale Learnable Graph Convolutional Networks 1416–1424.
- Gilmer, J.; Schoenholz, S. S.; Riley, P. F.; Vinyals, O.; and Dahl, G. E. 2017. Neural message passing for quantum chemistry. *arXiv preprint arXiv:1704.01212*.
- Grover, A.; and Leskovec, J. 2016. node2vec: Scalable feature learning for networks. In *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*, 855–864.
- Guo, S.; Lin, Y.; Feng, N.; Song, C.; and Wan, H. 2019. Attention based spatial-temporal graph convolutional networks for traffic flow forecasting. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, 922–929.
- Hajiramezani, E.; Hasanzadeh, A.; Narayanan, K.; Duffield, N.; Zhou, M.; and Qian, X. 2019. Variational graph recurrent neural networks. In *Advances in neural information processing systems*, 10701–10711.
- Hamilton, W.; Ying, Z.; and Leskovec, J. 2017. Inductive representation learning on large graphs. In *Advances in neural information processing systems*, 1024–1034.
- He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 770–778.
- Hu, W.; Fey, M.; Zitnik, M.; Dong, Y.; Ren, H.; Liu, B.; Catasta, M.; and Leskovec, J. 2020. Open graph benchmark: Datasets for machine learning on graphs. *arXiv preprint arXiv:2005.00687*.
- Ioffe, S.; and Szegedy, C. 2015. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*.
- Kingma, D. P.; and Ba, J. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Kipf, T. N.; and Welling, M. 2016a. Semi-supervised classification with graph convolutional networks. *ICLR*.
- Kipf, T. N.; and Welling, M. 2016b. Variational graph auto-encoders. *arXiv preprint arXiv:1611.07308*.
- Klicpera, J.; Bojchevski, A.; and Günnemann, S. 2019. Predict then propagate: Graph neural networks meet personalized pagerank. *ICLR*.
- Lee, J.; Lee, I.; and Kang, J. 2019. Self-Attention Graph Pooling. In *International Conference on Machine Learning*, 3734–3743.
- Li, C.; and Goldwasser, D. 2019. Encoding social information with graph convolutional networks for Political perspective detection in news media. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 2594–2604.
- Li, G.; Muller, M.; Thabet, A.; and Ghanem, B. 2019a. Deepgcns: Can gcns go as deep as cnns? In *Proceedings of the IEEE International Conference on Computer Vision*, 9267–9276.
- Li, G.; Xiong, C.; Thabet, A.; and Ghanem, B. 2020. DeeperGCN: All You Need to Train Deeper GCNs. *arXiv preprint arXiv:2006.07739*.
- Li, J.; Han, Z.; Cheng, H.; Su, J.; Wang, P.; Zhang, J.; and Pan, L. 2019b. Predicting Path Failure In Time-Evolving Graphs. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 1279–1289.
- Li, Q.; Han, Z.; and Wu, X.-M. 2018. Deeper Insights into Graph Convolutional Networks for Semi-Supervised Learning. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence (AAAI-18)*, 3538–3545. Association for the Advancement of Artificial Intelligence.
- Li, Y.; Tarlow, D.; Brockschmidt, M.; and Zemel, R. 2015. Gated graph sequence neural networks. *arXiv preprint arXiv:1511.05493*.
- Liu, M.; Gao, H.; and Ji, S. 2020. Towards Deeper Graph Neural Networks. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 338–348.
- Liu, Z.; Chen, C.; Li, L.; Zhou, J.; Li, X.; Song, L.; and Qi, Y. 2019. Geniepath: Graph neural networks with adaptive receptive paths. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, 4424–4431.
- Ma, Y.; Wang, S.; Aggarwal, C. C.; and Tang, J. 2019. Graph convolutional networks with eigenpooling. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 723–731.
- Monti, F.; Bronstein, M.; and Bresson, X. 2017. Geometric matrix completion with recurrent multi-graph neural networks. In *Advances in Neural Information Processing Systems*, 3697–3707.

- Oono, K.; and Suzuki, T. 2019. Graph neural networks exponentially lose expressive power for node classification. *arXiv preprint arXiv:1905.10947* .
- Page, L.; Brin, S.; Motwani, R.; and Winograd, T. 1999. The PageRank citation ranking: Bringing order to the web. Technical report, Stanford InfoLab.
- Paszke, A.; Gross, S.; Massa, F.; Lerer, A.; Bradbury, J.; Chanan, G.; Killeen, T.; Lin, Z.; Gimelshein, N.; Antiga, L.; et al. 2019. Pytorch: An imperative style, high-performance deep learning library. In *Advances in neural information processing systems*, 8026–8037.
- Perozzi, B.; Al-Rfou, R.; and Skiena, S. 2014. Deepwalk: Online learning of social representations. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, 701–710.
- Qiu, J.; Tang, J.; Ma, H.; Dong, Y.; Wang, K.; and Tang, J. 2018. Deepinf: Social influence prediction with deep learning. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2110–2119.
- Rong, Y.; Huang, W.; Xu, T.; and Huang, J. 2019. Dropedge: Towards deep graph convolutional networks on node classification. In *International Conference on Learning Representations*.
- Sen, P.; Namata, G.; Bilgic, M.; Getoor, L.; Galligher, B.; and Eliassi-Rad, T. 2008. Collective classification in network data. *AI magazine* 29(3): 93–93.
- Seo, Y.; Defferrard, M.; Vandergheynst, P.; and Bresson, X. 2018. Structured sequence modeling with graph convolutional recurrent networks. In *International Conference on Neural Information Processing*, 362–373. Springer.
- Tang, J.; Qu, M.; Wang, M.; Zhang, M.; Yan, J.; and Mei, Q. 2015. Line: Large-scale information network embedding. In *Proceedings of the 24th international conference on world wide web*, 1067–1077.
- Veličković, P.; Cucurull, G.; Casanova, A.; Romero, A.; Lio, P.; and Bengio, Y. 2018. Graph attention networks. *ICLR* .
- Wu, F.; Souza, A.; Zhang, T.; Fifty, C.; Yu, T.; and Weinberger, K. 2019. Simplifying Graph Convolutional Networks. In *International Conference on Machine Learning*, 6861–6871.
- Xhonneux, L.-P. A.; Qu, M.; and Tang, J. 2019. Continuous Graph Neural Networks. *arXiv preprint arXiv:1912.00967* .
- Xu, K.; Hu, W.; Leskovec, J.; and Jegelka, S. 2018a. How powerful are graph neural networks? *ICLR* .
- Xu, K.; Li, C.; Tian, Y.; Sonobe, T.; Kawarabayashi, K.-i.; and Jegelka, S. 2018b. Representation learning on graphs with jumping knowledge networks. *arXiv preprint arXiv:1806.03536* .
- Yan, T.; Zhang, H.; Li, Z.; and Xia, Y. 2020. Stochastic Graph Recurrent Neural Network. *arXiv preprint arXiv:2009.00538* .
- Yang, C.; Wang, R.; Yao, S.; Liu, S.; and Abdelzaher, T. 2020. Revisiting” Over-smoothing” in Deep GCNs. *arXiv preprint arXiv:2003.13663* .
- Yang, Z.; Cohen, W.; and Salakhudinov, R. 2016. Revisiting semi-supervised learning with graph embeddings. In *International conference on machine learning*, 40–48.
- Ying, R.; He, R.; Chen, K.; Eksombatchai, P.; Hamilton, W. L.; and Leskovec, J. 2018a. Graph convolutional neural networks for web-scale recommender systems. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 974–983.
- Ying, Z.; You, J.; Morris, C.; Ren, X.; Hamilton, W.; and Leskovec, J. 2018b. Hierarchical graph representation learning with differentiable pooling. In *Advances in neural information processing systems*, 4800–4810.
- Zhang, M.; and Chen, Y. 2017. Weisfeiler-lehman neural machine for link prediction. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 575–583.
- Zhang, M.; and Chen, Y. 2018. Link prediction based on graph neural networks. In *Advances in Neural Information Processing Systems*, 5165–5175.
- Zhang, M.; Cui, Z.; Neumann, M.; and Chen, Y. 2018. An end-to-end deep learning architecture for graph classification. In *Thirty-Second AAAI Conference on Artificial Intelligence*.
- Zhao, L.; and Akoglu, L. 2019. Pairnorm: Tackling over-smoothing in gnns. *arXiv preprint arXiv:1909.12223* .
- Zitnik, M.; and Leskovec, J. 2017. Predicting multicellular function through multi-layer tissue networks. *Bioinformatics* 33(14): i190–i198.