Graph Optimized Convolutional Networks

Bo Jiang, Ziyan Zhang, Jin Tang and Bin Luo

School of Computer Science and Technology, Anhui University jiangbo@ahu.edu.cn

Abstract

Graph Convolutional Networks (GCNs) have been widely studied for graph data representation and learning tasks. Existing GCNs generally use a fixed single graph which may lead to weak suboptimal for data representation/learning and are also hard to deal with multiple graphs. To address these issues, we propose a novel Graph Optimized Convolutional Network (GOCN) for graph data representation and learning. Our GOCN is motivated based on our re-interpretation of graph convolution from a regularization/optimization framework. The core idea of GOCN is to formulate graph optimization and graph convolutional representation into a unified framework and thus conducts both of them cooperatively to boost their respective performance in GCN learning scheme. Moreover, based on the proposed unified graph optimization-convolution framework, we propose a novel Multiple Graph Optimized Convolutional Network (M-GOCN) to naturally address the data with multiple graphs. Experimental results demonstrate the effectiveness and benefit of the proposed GOCN and M-GOCN.

1 Introduction

Convolutional Neural Networks (CNNs) have been widely applied for grid-like structure data representation and learning in computer vision and machine learning area. However, in many real applications, data are not coming with grid-like structure but instead have some irregular structures which are usually represented as structured graphs. Traditional CNNs generally fail to address graph-structured data.

Recently, Graph Convolutional Networks (GCNs) have been widely studied to deal with arbitrary graph-structured data representation and learning [4, 15, 19, 9, 30, 24, 33]. The aim of GCNs is trying to define some reasonable convolution operations on arbitrary structured graphs. For example, Bruna et al. [3] propose to define a graph convolution by using the eigen-decomposition of graph Laplacian matrix. Henaff et al. [10] further introduce a spatially constrained spectral filters to define graph convolution. Kipf et al. [15] propose to explore the first-order approximation of spectral filters and present a simple Graph Convolutional Network (GCN) for semi-supervised learning. Li et al. [30] present an adaptive graph CNNs, in which the graph is learned adaptively by employing a metric learning method. Hamilton et al. [9] present a general inductive representation and learning framework to generate embeddings for the unseen nodes. Velickovic et al. [30] propose Graph Attention Networks (GAT) for graph based semi-supervised learning. Klicpera et al. [13] propose to combine GCN and PageRank together to derive an improved propagation scheme in layer-wise propagation. Some recent works also explore specific graph neural networks for computer vision tasks [18, 28, 12, 29].

The above GCNs have been widely used for graph data representation and learning. One important aspect of GCNs is the graph structure representation of data. In general, the graph data we feed to existing GCNs should be a single graph which is obtained from either domain knowledge (e.g., social network) or human establishment, such as k-NN graph. However, there are three main issues. First, traditional human established graphs generally use fixed parameters to determine the graph

structure and thus are usually sensitive to the local noises and errors. Second, the graphs obtained from domain knowledge or established by human are generally independent of graph convolutional learning, which thus are not guaranteed to be optimal for graph convolutional representation/learning in GCNs. Third, existing GCNs are generally hard to deal with multiple graphs, although some heuristic fusion strategies can be utilized for multiple graph convolutional network learning [27, 25, 35]. It is known that the convolution operation on multiple graphs is not as well-defined as on single graph.

To address these issues, in this paper, we propose a novel Graph Optimized Convolutional Network (GOCN) for data representation and learning problem. Our GOCN is motivated based on our new re-interpretation of graph convolution from a regularization/optimization framework. The core idea of GOCN is to formulate graph optimization and graph convolutional representation into a unified framework and thus conducts both of them cooperatively to boost their respective performance in GCN learning scheme. The main advantage of GOCN is that the learned representation of data can provide useful "weakly" supervised information for learning a better graph which simultaneously facilitates graph convolutional representation and learning. Furthermore, based on the proposed unified graph optimization-convolution framework, we propose a novel Multiple Graph Optimized Convolutional Network (M-GOCN) to naturally address the data with multiple graphs.

Overall, the main contributions of this paper are summarized as follows:

- We propose to reformulate graph convolution learning as a regularization framework, based
 on which a unified graph optimization-convolution (GOC) framework is derived to learn
 an optimal graph for graph convolutional representation and learning.
- Based on the proposed unified GOC model, we propose a novel Graph Optimized Convolutional Network (GOCN) which conducts both graph construction and graph convolution simultaneously in GCN scheme for data representation and semi-supervised learning.
- We extend GOC to deal with multiple graphs and provide a novel Multiple Graph Optimized Convolutional Network (M-GOCN) for multi-graph data representation and learning.

Experimental results on several datasets demonstrate the effectiveness of the proposed GOCN and M-GOCN methods.

2 Revisiting GCN

Recently, Graph Convolutional Networks (GCNs) have been widely studied for graph data representation and learning [4, 15, 10, 30]. The core aspect of GCNs is the specific definition of graph convolution in layer-wise propagation. Here, we briefly review the widely used GCN model proposed in work [15]. Given an input feature matrix $H^{(0)} = X \in \mathbb{R}^{n \times d_0}$ and graph $A \in \mathbb{R}^{n \times n}$ with $A_{ii} = 0$, GCN defines the layer-wise propagation as [15],

$$H^{(k+1)} = \sigma((I + D^{-\frac{1}{2}}AD^{-\frac{1}{2}})H^{(k)}\Theta^{(k)})$$
(1)

where $k=0,1,\cdots K-1$ and I denotes an identity matrix, and D is a diagonal degree matrix. $\Theta^{(k)} \in \mathbb{R}^{d_k \times d_{k+1}}$ is a layer-specific trainable weight matrix, and $\sigma(\cdot)$ denotes an activation function.

The last layer of GCN outputs the final representation $H^{(K)}$ of graph nodes, which can be used for many learning tasks, such as clustering, visualization and (semi-supervised) classification etc. In this paper, we focus on semi-supervised classification. For this task, a softmax activation function is further used to output the label prediction $P \in \mathbb{R}^{n \times c}$ for graph nodes, where c denotes class number. The weight matrices of GCN network $\{\Theta^{(0)},\Theta^{(1)},\cdots\Theta^{(K-1)}\}$ are optimized by minimizing the cross-entropy loss as [15],

$$\mathcal{L}_{\text{Semi-GCN}} = -\sum_{i \in L} \sum_{j=1}^{c} Y_{ij} \ln P_{ij}$$
 (2)

where L indicates the set of labeled nodes and each row $Y_{i\cdot}$, $i \in L$ of Y denotes the corresponding label indication vector for the i-th labeled node.

3 Graph Optimized Convolutional Network

In this section, we present our Graph Optimized Convolutional Network (GOCN) model. GOCN is motivated based on our re-interpretation of GCN by using a regularization framework. In the following, we first present our regularization reformulation of graph convolution in §3.1. Based on it, we then derive a unified framework of graph optimization-convolution operation in §3.2. Finally, we present our GOCN architecture in §3.3.

3.1 Regularization framework of GCN

The propagation rule Eq.(1) in GCN can be decomposed into two operations, i.e.,

$$Z^{(k)} = (I + D^{-\frac{1}{2}}AD^{-\frac{1}{2}})H^{(k)}$$
(3)

$$H^{(k+1)} = \sigma(Z^{(k)}\Theta^{(k)}) \tag{4}$$

where Eq.(3) defines a kind of *feature aggregation* for features $H^{(k)}$ on graph and Eq.(4) gives a non-linear feature transformation via projection $\Theta^{(k)}$ and non-linear activation $\sigma(\cdot)$. For simplicity, we rewrite Eq.(3) as

$$Z = (I + \hat{A})H = \hat{A}H + H \tag{5}$$

where $\hat{A} = D^{-\frac{1}{2}}AD^{-\frac{1}{2}}$ and $\hat{A}_{ii} = 0$ (because $A_{ii} = 0$). From Eq.(5), we can note that GCN employs an one-step feature aggregation on normalized graph \hat{A} (biased by feature itself) to obtain contextual feature representation in layer-wise propagation.

First, here one can also explore a more flexible T-step aggregation in GCN layer-wise propagation as

$$Z^{(t+1)} = \alpha \hat{A} Z^{(t)} + (1 - \alpha) H = \left[(\alpha \hat{A})^{t+1} + (1 - \alpha) \sum_{i=0}^{t} (\alpha \hat{A})^{i} \right] H$$
 (6)

where $t=0,1\cdots T-1$ and $Z^{(0)}=H$. Parameter $\alpha\in(0,1)$ denotes the fraction of feature information that node v_i receives from its neighbors on normalized graph \hat{A} . Note that, when t=0 and $\alpha=0.5$, the aggregation (Eq.(6)) degenerates to GCN update Eq.(5). When $t\to\infty$, this aggregation converges to an equilibrium solution as [6,5],

$$Z^* = (1 - \alpha)(I - \alpha \hat{A})^{-1}H \tag{7}$$

Second, one important property of the update Eq.(6) is that it can be theoretically explained by using an optimization framework [6]. Specifically, the converged solution of Eq.(7) is the optimal solution that minimizes the following optimization problem,

$$\min_{Z} \mathcal{R}(\hat{A}, H; Z) = \text{Tr}(Z^{T}(I - \hat{A})Z) + \mu \|Z - H\|_{F}^{2}$$
(8)

where $\mu=\frac{1}{\alpha}-1$ is a replacement parameter of α to balance two terms. $\mathrm{Tr}(\cdot)$ denotes the trace function and $\|\cdot\|_F$ denotes the Frobenious norm. Moreover, the update Eq.(6) (or Eq.(7)) also provides an approximate solution to this problem by using a T-step power iteration algorithm.

In particular, we can note that GCN update (Eq.(5)) provides a very approximate solution for the problem $\mathcal{R}(\hat{A}, H; Z)$ by using an one-step power algorithm. This provides a regularization/optimization framework interpretation for GCN update rule, which motivates us to develop some more effective graph convolution (feature aggregation) variants in GCN layer-wise propagation. In this paper, we focus on graph construction \hat{A} and aim to learn a more effective graph for GCN representation/learning. To do so, in the following we first derive a unified framework of graph optimization and convolution, followed by a simple power iteration implementation. Based on this unified framework, we then develop our GOCN architecture in §3.3.

3.2 Unified graph optimization-convolution model

One important aspect of the above feature aggregation is the graph construction \hat{A} . Constructing a good graph to represent data relationship is generally important for data representation and (semi-supervised) learning tasks. Existing GCNs generally use a fixed graph which may lead to weak

suboptimal for data representation and learning. Our aim in this section is to propose a unified graph optimization-convolution model that aims to learn an optimal graph adaptively for feature aggregation in GCN scheme.

Formally, given an initial (normalized) graph \hat{A} , we propose to learn an optimal graph S that best serves the above feature aggregation Eq.(8) while preserving the structure information encoded in \hat{A} . This can be achieved by optimizing the following unified framework,

$$\min_{S,Z} \mathcal{U}_{goc} = \mathcal{G}(\hat{A}; S) + \gamma \mathcal{R}(S, H; Z)$$

$$s.t. \quad S = S^T, S_{ij} \ge 0$$
(9)

where $\mathcal{G}(\hat{A}; S)$ denotes the graph learning functions and γ is used to balance two terms. In this paper, we use the simple Frobenious norm and propose a graph optimized feature aggregation as,

$$\min_{S,Z} \mathcal{U}_{goc} = \|\hat{A} - S\|_F^2 + \gamma \mathcal{R}(S, H; Z)$$

$$s.t. \quad S = S^T, S_{ij} \ge 0$$

$$(10)$$

The optimal S and Z can be obtained via a simple algorithm which alternatively conducts the following step 1 and step 2 until convergence.

Step 1. Solving S while fixing Z, the problem becomes

$$\min_{S} \|\hat{A} - S\|_F^2 + \gamma \text{Tr}(Z^T (I - S)Z) \tag{11}$$

$$s.t. \ S = S^T, S_{ij} \ge 0$$

It has a simple closed-form solution which is given as

$$S_{ij} = \max \left\{ (\hat{A} + \frac{\gamma}{2} Z Z^T)_{ij}, 0 \right\}$$
 (12)

Step 2. Solving Z while fixing S, the problem becomes to Eq.(8). The optimal solution is

$$Z = (1 - \alpha)(I - \alpha S)^{-1}H \tag{13}$$

and an approximate solution can be obtained by

$$Z = [(\alpha S)^{T} + (1 - \alpha) \sum_{i=0}^{T-1} (\alpha S)^{i}] H$$
(14)

Remark. Here, to avoid the inversion computation in Eq.(13), we can instead use a T-step power iteration and obtain Eq.(14) to compute the optimal Z approximately [6].

3.3 GOCN architecture

Using the proposed unified model \mathcal{U}_{goc} , we present our Graph Optimized Convolutional Network (GOCN). Similar to the architecture of standard GCN [15], our GOCN contains one input layer, several hidden propagation layers and one final perceptron layer, as introduced in the following.

3.3.1 Hidden Propagation Layer

For hidden propagation layer, it takes features $H^{(k)} \in \mathbb{R}^{n \times d_k}$ and an initial graph $\hat{A} \in \mathbb{R}^{n \times n}$ as the input and outputs feature map $H^{(k+1)} \in \mathbb{R}^{n \times d_{k+1}}$. Let $Z^* = \Phi(A, H)$ be the optimal Z-solution of unified model \mathcal{U}_{qoc} (Eq.(10)), then GOCN conducts layer-wise propagation as

$$H^{(k+1)} = \sigma(\Phi(\hat{A}, H^{(k)})\Theta^{(k)})$$
(15)

where $k=0,1,\cdots K-1$ and $\Theta^{(k)}\in\mathbb{R}^{d_k\times d_{k+1}}$ denotes the layer-specific trainable weight matrix. $\sigma(\cdot)$ denotes an activation function.

Efficient computation of $\Phi(\hat{A}, H^{(k)})$. Exactly calculating $\Phi(\hat{A}, H^{(k)})$ is time consuming due to (i) inversion operation (Eq.(13)) and (ii) alternative computation of step 1 and step 2 until convergence. To overcome this issue, we first use update Eq.(14) to compute the optimal Z in Eq.(13) approximately. Second, we use a M-step alternative iteration to optimize S and S approximately. The detail propagation rule in GOCN hidden layer is summarized in Algorithm 1. In our experiments, we set S and S are included in the constant S and S and S are included in the constant S and S

Algorithm 1 GOCN layer-wise propagation

- 1: **Input:** Feature matrix $H^{(k)} \in \mathbb{R}^{n \times d_k}$, initial graph $\hat{A} \in \mathbb{R}^{n \times n}$ and weight parameter $\Theta^{(k)}$, parameters γ, μ , maximum iteration T, M
- 2: Output: Feature map $H^{(k+1)} \in \mathbb{R}^{n \times d_{k+1}}$
- 3: Initialize $Z = H^{(k)}$
- 4: **for** $t_m = 1, 2 \cdots M$ **do**
- 5: Compute S as

$$S_{ij} = \max\left\{ \left(\hat{A} + \frac{\gamma}{2} Z Z^T \right)_{ij}, 0 \right\}$$

6: Compute Z as

$$Z = \left[(\alpha S)^T + (1 - \alpha) \sum_{i=0}^{T-1} (\alpha S)^i \right] H^{(k)}$$

- 7. end for
- 8: Return $H^{(k+1)} = \sigma(Z\Theta^{(k)})$

3.3.2 Final Perceptron Layer

The last layer of GOCN outputs the final label prediction $P \in \mathbb{R}^{n \times c}$ for graph nodes, where c denotes the number of class. Similar to work [15], the optimal network weight parameters $\{\Theta^{(0)}, \Theta^{(1)}, \cdots \Theta^{(K-1)}\}$ of GOCN are obtained by minimizing the following cross-entropy loss function over all the labeled nodes L, i.e.,

$$\mathcal{L}_{\text{Semi-GOCN}} = -\sum_{i \in L} \sum_{j=1}^{c} Y_{ij} \ln P_{ij}$$
 (16)

where L indicates the set of labeled nodes and $Y_{i\cdot}$, $i\in L$ denotes the corresponding label indication vector for the i-th labeled node.

4 GOCN on Multiple Graphs

One desired property of the proposed GOCN is that it can be naturally adapted to address the data with multiple graph representations. This is because we can extend the unified model \mathcal{U}_{goc} (Eqs.(9,10)) to deal with multiple graphs.

4.1 Multi-graph convolution model

Given an input feature matrix $X = H^{(0)} \in \mathbb{R}^{n \times d_0}$ with multiple graphs $\mathcal{A} = \{\hat{A}^{(1)}, \hat{A}^{(2)} \cdots \hat{A}^{(m)}\}$, we can obtain an optimal feature aggregation $Z = \Phi_m(\mathcal{A}, H)$ on graph set \mathcal{A} by optimizing,

$$\min_{S,Z} \mathcal{U}_{mgoc} = \mathcal{G}(\mathcal{A}; S) + \gamma \mathcal{R}(S, H; Z)$$
s.t. $S = S^T, S_{ij} \ge 0$ (17)

In particular, we use Frobenious norm and propose our multiple graph optimized feature aggregation as

$$\min_{S,Z,w} \mathcal{U}_{mgoc} = \sum_{v=1}^{m} w_v^r \|\hat{A}^{(v)} - S\|_F^2 + \gamma \mathcal{R}(S, H; Z)$$

$$s.t. \quad \sum_{v=1}^{m} w_v = 1, w_v \ge 0, S = S^T, S_{ij} \ge 0$$
(18)

where $w=(w_1,w_2\cdots w_m)$ denote the important weights of different graphs which are learned adaptively. The parameter r>1 is used to control the weight distribution, as suggested in previous work [32]. The optimal S, Z and w can be obtained by alternatively conducting the following Step $1\sim3$ until convergence.

Step 1. Solving S while fixing Z, w, the problem becomes

$$\min_{S} \sum_{v=1}^{m} w_{v}^{r} \|\hat{A}^{(v)} - S\|_{F}^{2} + \gamma \text{Tr}(Z^{T}(I - S)Z)$$
s.t. $S = S^{T}, S_{ij} \ge 0$ (19)

It has a simple closed-form solution as

$$S_{ij} = \max\left\{ \left(\sum_{v=1}^{m} w_v^r \hat{A}^{(v)} + \frac{\gamma}{2} Z Z^T \right)_{ij}, 0 \right\}$$
 (20)

Step 2. Solving Z while fixing S, w, the problem becomes to Eq.(8). The optimal solution is given as Eq.(13) exactly or Eq.(14) approximately.

Step 3. Solving w while fixing S, Z, the problem becomes

$$\min_{w} \sum_{v=1}^{m} w_{v}^{r} \|\hat{A}^{(v)} - S\|_{F}^{2} + \gamma \mathcal{R}(S, H; Z)
s.t. \sum_{v=1}^{m} w_{v} = 1, w_{v} \ge 0$$
(21)

The Lagrangian function is

$$\min_{w} \sum_{v=1}^{m} w_{v}^{r} \|\hat{A}^{(v)} - S\|_{F}^{2} + \xi \left(\sum_{v=1}^{m} w_{v} - 1\right)$$

$$s.t. \ w_{v} > 0$$
(22)

where ξ is the Lagrangian multiplier. With some simple algebraic manipulations, the optimal w is derived as

$$w_v = \frac{(1/\|\hat{A}^{(v)} - S\|_F^2)^{1/(r-1)}}{\sum_{v=1}^m (1/\|\hat{A}^{(v)} - S\|_F^2)^{1/(r-1)}}$$
(23)

M-GOCN architecture

Let $\mathcal{A}=\{\hat{A}^{(1)},\hat{A}^{(2)}\cdots\hat{A}^{(m)}\}$ and $Z=\Phi_m(\mathcal{A},H)$ be the optimal solution of problem Eq.(18), then our Multiple GOCN (M-GOCN) conducts the layer-wise propagation as,

$$H^{(k+1)} = \sigma(\Phi_m(\mathcal{A}, H^{(k)})\Theta^{(k)})$$
(24)

where $k = 0, 1 \cdots K - 1$ and $\Theta^{(k)}$ denotes the layer-specific trainable weight matrix. $\sigma(\cdot)$ denotes an activation function, such as $ReLU(\cdot) = max(0, \cdot)$.

Remark. Similar to GOCN, here we can instead use a T-step power iteration to compute the optimal Z approximately to avoid the inversion computation. The complete algorithm to compute $\Phi_m(\mathcal{A}, H^{(k)})$ in M-GOCN hidden propagation is summarized in Algorithm 2. In our experiments, we set T=2 and M=3 for M-GOCN.

Algorithm 2 M-GOCN layer-wise propagation

- 1: Input: Feature matrix $H^{(k)} \in \mathbb{R}^{n \times d_k}$, initial multiple graphs $\mathcal{A} = \{\hat{A}^{(1)}, \hat{A}^{(2)} \cdots \hat{A}^{(m)}\}$ and weight parameter $\Theta^{(k)} \in \mathbb{R}^{n \times d_{k+1}}$, parameters γ, μ, r , maximum iteration T, M
- 2: **Output:** Feature map $H^{(k+1)} \in \mathbb{R}^{n \times d_{k+1}}$
- 3: Initialize $Z = H^{(k)}$; $w = (1/m \cdots 1/m)$
- 4: **for** $t_m = 1, 2 \cdots M$ **do**
- Compute S as

Compute
$$S$$
 as
$$S_{ij} = \max \left\{ \left(\sum_{v=1}^{m} w_v^r \hat{A}^{(v)} + \frac{\gamma}{2} Z Z^T \right)_{ij}, 0 \right\}$$
Compute Z as

$$Z = \left[(\alpha S)^T + (1 - \alpha) \sum_{i=0}^{T-1} (\alpha S)^i \right] H^{(k)}$$

7:

$$w_v = \frac{(1/\|\hat{A}^{(v)} - S\|_F^2)^{1/(r-1)}}{\sum_{v=1}^{m} (1/\|\hat{A}^{(v)} - S\|_F^2)^{1/(r-1)}}$$

- 9: Return $H^{(k+1)} = \sigma(Z\Theta^{(k)})$

The perceptron layer of M-GOCN outputs the final label prediction $P \in \mathbb{R}^{n \times c}$ where c denotes the number of class, and the loss function is designed as a cross-entropy loss defined over labeled data which is the same as GOCN, as discussed in §3.3.2.

5 Experiments

In order to evaluate the effectiveness of the proposed GOCN and M-GOCN, we test them on several datasets and compare them with some other baseline models.

5.1 Datasets

In single graph learning experiments, we test GOCN on six datasets including three standard network datasets (Citeseer, Cora and Pubmed [26]) and three image datasets (CIFAR10 [16], SVHN [20] and Scene 15 [11]). Their usages in our experiments are introduced below.

Citeseer contains 3327 nodes and 4732 edges whose nodes denote documents and edges encode the citation relationships between documents. Each node is represented by a 3703 dimension feature descriptor and all nodes are classified into six classes.

Cora contains 2708 nodes and 5429 edges. Each node is represented by a 1433 dimension feature descriptor and all the nodes are falling into six classes.

Pubmed contains 19717 nodes and 44338 edges. Each node is represented by a 500 dimension feature descriptor and all the nodes are falling into three classes.

Scene15 dataset consists of 4485 scene images with 15 different categories [11]. For each image, we use the feature descriptor provided by work [11].

CIFAR10 dataset contains 50000 natural RGB color images and all images are falling into 10 classes [16]. In our experiments, we select 10000 images in all with 1000 images per class for evaluation. We have not use all of images because large storage and high computational complexity are required for graph convolution operation in our GOCN and other comparing GCN based methods. For each image, a CNN feature descriptor is extracted to represent it.

SVHN dataset contains 73257 training and 26032 test RGB images [20]. Similarly, we select 10000 images in all with 1000 images for each class in our evaluation. For each image, a CNN feature descriptor is extracted.

In our multi-graph learning experiments, we evaluate M-GOCN on three datasets including MSRC-v1 [31], Caltech101-7 [17, 22] and Handwritten numerals [1]. Their usages in our experiments are introduced below.

MSRC-v1 [31] contains 8 classes of 240 images. Each class contains 30 images. Similar to the setup in work [22], we obtain five graphs by using five different kinds of visual descriptors and the final input feature of each graph node is constructed by concatenating different descriptors together. Caltech101-7 [17] contains 101 categories of images. Following the setup of work [21], we select the widely used 7 classes and obtain 1474 images in all in our experiments. We construct six graphs by using six different feature descriptors and obtain the input feature of each graph node by concatenating these descriptors together [21].

Handwritten numerals [1] contains 2000 digits from '0' to '9' and each digit class has 200 data samples. We construct six graphs by using six published feature descriptors and obtain the final input feature of each graph node by concatenating them together.

5.2 Experimental setup

5.2.1 Parameter setting

Similar to experimental setting of GCN [15], we use a two-layer graph convolutional network and the number of units in hidden layer is set to 16. We train our GOCN using an ADAM algorithm [14] with 10000 maximum epochs and learning rate of 0.01. We stop training if the validation loss does not decrease for 100 consecutive epochs, as suggested in [15]. All the network weights are initialized using Glorot initialization [8]. The parameters α , γ are set to 0.9 and 20 respectively. Note that, GOCN is not insensitive to these parameters. We provide additional experiments across different settings of α , β and hidden layer number in §5.4.

5.2.2 Data setting

For network datasets (Citeseer, Cora and Pubmed), we utilize the similar data setting used in previous works [15, 30]. For each class, we select 20 nodes as labeled data and 300 nodes as validation data, and then evaluate the performance of label prediction on the remaining 1000 test nodes. For

image dataset CIFAR10 and SVHN, we randomly select 1000, 2000 and 3000 images as labeled samples and use the remaining data as unlabeled samples. For unlabeled samples, we select 1000 images for validation purpose and use the remaining 8000, 7000 and 6000 images as test samples, respectively. For image dataset Scene15 [11], we randomly select 500, 750 and 1000 images as labeled data and select 500 images for validation. The remaining images are used for testing. All the reported results are averaged over five runs with different data splits of training, validation and testing. In our multi-graph learning experiments, for all datasets (MSRC-v1, Caltech101-7 and Handwritten numerals), we select 10%, 20% and 30% nodes as labeled samples and use the remaining data as unlabeled samples. For unlabeled samples, we also use 5% nodes for validation purpose to determine the convergence criterion, and use the remaining 85%, 75% and 65% nodes respectively as test samples. All the reported results are averaged over five runs with different data splits of training, validation and testing.

Table 1: Comparison results of different methods on dataset Citeseer, Cora and Pubmed. The best results are marked by bold.

Methond	Citeseer	Cora	Pubmed
ManiReg	60.1%	59.5%	70.7%
LP	59.6%	59.0%	71.1%
DeepWalk	43.2%	67.2%	65.3%
DGI	71.5%	76.8%	77.2%
GCN	68.9%	82.9%	77.9%
GAT	71.0%	83.2%	78.0%
GOCN	71.8%	84.8%	79.7%

5.3 Comparison with state-of-the-art methods

5.3.1 Evaluation on single graph

We first compare our GOCN with the baseline model GCN [15] to demonstrate the benefit of graph optimization. Also, we compare our method against some other graph neural network based semisupervised learning methods which contain two traditional graph based semi-supervised learning methods including Label Propagation (LP) [34], Manifold Regularization (ManiReg) [2], and four graph neural network methods including DeepWalk [23], Graph Convolutional Network (GCN) [15], Graph Attention Networks (GAT) [30] and Deep Graph Informax (DGI) [24]. Table 1 shows the comparison results on three citation network benchmark datasets. Table 2 and 3 summarize the comparison results on three image datasets. The best results are marked as bold. We can note that (1) Comparing with the baseline model GCN [15], GOCN obtains obviously better learning results on all datasets, especially on image datasets. This demonstrates the higher predictive ability of GOCN on semi-supervised classification by incorporating graph optimization, which indicates that GOCN conducts data representation and semi-supervised learning more optimal than GCN. (2) GOCN performs better than recent graph network GAT [30] and DGI [24], which demonstrates the benefit of GOCN on data representation and learning. (3) GOCN generally performs better than other graph based semi-supervised method LP [34], ManiReg [2] and DeepWalk [23], which further indicates the effectiveness of GOCN on conducting semi-supervised classification tasks.

Table 2: Comparison results of different methods on Scene15 dataset. The best results are marked by bold.

Dataset	Scene15				
No. of label	500	750	1000		
ManiReg	81.29 ± 3.35	86.45 ± 1.91	89.86 ± 0.71		
LP	89.40 ± 4.74	92.12 ± 2.87	92.98 ± 2.45		
Deep Walk	95.64 ± 0.24	96.01 ± 0.24	96.53 ± 0.37		
ĎGI	92.94 ± 1.61	94.21 ± 0.64	94.27 ± 0.94		
GCN	91.42 ± 2.07	94.41 ± 0.92	95.44 ± 0.89		
GAT	93.98 ± 0.75	94.64 ± 0.41	95.03 ± 0.46		
GOCN	95.87 ± 0.56	97.40±0.34	98.00±0.29		

Table 3: Comparison results of different methods on dataset SVHN and CIFAR10. The best results are marked by bold.

Dataset	SVHN			CIFAR10		
No. of label	1000	2000	3000	1000	2000	3000
ManiReg	69.44±0.69	72.73 ± 0.44	74.63 ± 0.45	52.30 ± 0.66	57.08 ± 0.80	59.69 ± 0.71
LP	69.68 ± 0.84	70.35 ± 1.73	69.47 ± 2.96	57.52 ± 0.67	59.22 ± 0.67	60.38 ± 0.51
Deep Walk	74.64 ± 0.23	76.21 ± 0.23	77.04 ± 0.42	56.16 ± 0.54	59.73 ± 0.35	61.26 ± 0.32
DGI	70.82 ± 1.22	72.83 ± 0.79	73.16 ± 1.20	58.97 ± 0.61	60.26 ± 0.56	60.56 ± 0.36
GCN	71.33 ± 1.48	73.43 ± 0.46	73.63 ± 0.52	60.43 ± 0.56	60.91 ± 0.50	60.99 ± 0.49
GAT	73.87 ± 0.32	74.85 ± 0.55	75.17 ± 0.43	63.25 ± 0.50	65.55 ± 0.58	66.56 ± 0.58
GOCN	80.72±0.35	82.67±0.25	83.63±0.37	68.13±0.58	71.83 ± 0.37	73.66 ± 0.52

Table 4: Comparison results of different multi-graph learning methods on dataset Caltech101-7, MSRC-v1 and Handwritten numerals, respectively. The best results are marked by bold.

Dataset	Caltech101-7		MSRC-v1		Handwritten numerals	
Ratio of label	10%	20%	10%	20%	10%	20%
GCN(1)	82.62±0.78	84.71±1.38	62.08 ± 6.74	66.21±5.42	90.92 ± 0.22	91.33±0.45
GCN(2)	85.00 ± 1.90	86.34 ± 1.73	81.97 ± 1.61	$85.84{\pm}2.90$	94.72 ± 0.22	95.39 ± 0.77
GCN(3)	86.99 ± 1.42	88.37 ± 0.82	84.94 ± 3.36	88.07 ± 1.54	95.59 ± 0.61	96.16 ± 0.57
GCN(4)	93.18 ± 0.98	93.44 ± 0.55	77.47 ± 3.20	83.10 ± 1.54	95.78 ± 0.31	96.45 ± 0.68
GCN(5)	92.35 ± 0.62	92.49 ± 0.47	80.00 ± 2.60	84.72 ± 3.46	88.53 ± 0.53	89.17 ± 0.63
GCN(6)	92.05 ± 0.78	92.45 ± 0.86	-	-	82.34 ± 0.67	83.11 ± 0.25
GCN-M	90.39 ± 2.02	91.91 ± 0.33	86.81 ± 3.60	90.31 ± 1.87	96.42 ± 0.47	97.08 ± 0.48
Multi-GCN	95.09 ± 0.62	96.08 ± 0.58	87.14 ± 2.13	90.43 ± 1.45	97.14 ± 0.23	97.88 ± 0.30
MLAN	93.45±0.36	94.86±0.29	83.42±2.10	87.27 ± 1.66	97.23 ± 0.26	97.46 ± 0.52
M-GOCN	95.97±0.36	97.23±0.47	90.22±1.06	91.68±1.50	97.76±0.42	97.92±0.41

5.3.2 Evaluation on multiple graphs

For multi-graph learning tasks, we compare our M-GOCN against the state-of-the art baseline methods which contain GCN(v) that conducts traditional GCN [15] on the v-th singe graph $\hat{A}^{(v)}$ and node content features X; GCN-M that conducts traditional GCN [15] on the averaged graph representation $\bar{A} = \frac{1}{m} \sum_{v=1}^{m} \hat{A}^{(v)}$ and node features X; Multi-GCN that first learns representations for multiple graphs $\hat{A}^{(v)}$ by using/sharing the common parameters (as suggested in work [7]), and then select the final representation with the lowest training loss function for multi-graph representation; MLAN [21] that is a multi-view learning model for graph learning and semi-supervised classification. Table 4 summarizes the comparison results of semi-supervised classification on these datasets. Here, we can note that (1) The proposed M-GOCN performs obviously better than traditional GCN [15] model conduced on each individual graph $A^{(v)}$. This clearly demonstrates the effectiveness of the proposed M-GOCN on learning a compact representation for multiple graphs by intergrading the information of multiple graphs together. (2) M-GOCN performs better than the baseline method GCN-M, Multi-GCN [7] and MLAN [21], which indicates the effectiveness of M-GOCN architecture on conduct multiple graph representation and semi-supervised learning tasks.

6 Conclusion and Future Works

This paper proposes a novel Graph Optimized Convolutional Network (GOCN) for graph data representation and semi-supervised learning. GOCN is inspired based on the re-formulation of graph convolution as a regularization framework. GOCN integrates graph optimization and convolution in a unified scheme and thus can boost their respectively performance in graph neural network learning. Also, GOCN can be naturally extended to M-GOCN to deal with multiple graphs. Experimental results on several benchmarks demonstrate the effectiveness and benefits of the proposed GOCN and M-GOCN on semi-supervised learning tasks. In the future, we will explore GOCN and M-GOCN methods on some other learning tasks, such as graph clustering, embedding etc. Also, we will adapt them on some more computer vision tasks, such as object detection, image co-segmentation and multiple object tracking etc.

References

- [1] D. N. Arthur Asuncion. Uci machine learning repository. 2007.
- [2] M. Belkin, P. Niyogi, and V. Sindhwani. Manifold regularization: A geometric framework for learning from labeled and unlabeled examples. *Journal of machine learning research*, 7(Nov):2399–2434, 2006.
- [3] J. Bruna, W. Zaremba, A. Szlam, and Y. LeCun. Spectral networks and locally connected networks on graphs. In *International Conference on Learning Representations*, 2014.
- [4] M. Defferrard, X. Bresson, and P. Vandergheynst. Convolutional neural networks on graphs with fast localized spectral filtering. In Advances in Neural Information Processing Systems, pages 3844–3852, 2016
- [5] A. G. O. B. Dengyong Zhou, Jason Weston and B. Schölkopf. Ranking on data manifolds. In NIPS, pages 169–176, 2004.
- [6] T. N. L. J. W. Dengyong Zhou, Olivier Bousquet and B. Schölkopf. Learning with local and global consistency. In NIPS, pages 321–328, 2004.
- [7] D. K. Duvenaud, D. Maclaurin, J. Iparraguirre, R. Bombarell, T. Hirzel, A. Aspuru-Guzik, and R. P. Adams. Convolutional networks on graphs for learning molecular fingerprints. In *Advances in Neural Information Processing Systems*, pages 2224–2232, 2015.
- [8] X. Glorot and Y. Bengio. Understanding the difficulty of training deep feedforward neural networks. In *International conference on artificial intelligence and statistics*, pages 249–256, 2010.
- [9] W. Hamilton, Z. Ying, and J. Leskovec. Inductive representation learning on large graphs. In *Advances in Neural Information Processing Systems*, pages 1024–1034, 2017.
- [10] M. Henaff, J. Bruna, and Y. LeCun. Deep convolutional networks on graph-structured data. *arXiv preprint* arXiv:1506.05163, 2015.
- [11] Z. Jiang, Z. Lin, and L. S. Davis. Label consistent k-svd: Learning a discriminative dictionary for recognition. IEEE transactions on pattern analysis and machine intelligence, 35(11):2651–2664, 2013.
- [12] L. W. Y. W. J. D. Jiayuan Gu, Han Hu. Learning region features for object detection. In *European Conference on Computer Vision (ECCV)*, pages 381–395, 2018.
- [13] A. B. Johannes Klicpera and S. Günnemann. Combining neural networks with personalized pagerank for classification on graphs. In *ICLR*, 2019.
- [14] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015.
- [15] T. N. Kipf and M. Welling. Semi-supervised classification with graph convolutional networks. *arXiv* preprint arXiv:1609.02907, 2016.
- [16] A. Krizhevsky and G. Hinton. Learning multiple layers of features from tiny images. Technical report, Citeseer, 2009.
- [17] Y. Li, F. Nie, H. Huang, and J. Huang. Large-scale multi-view spectral clustering via bipartite graph. In *AAAI*, pages 2750–2756, 2015.
- [18] D.-A. H. S. S. S. Y. Michelle Guo, Edward Chou and L. Fei-Fei. Neural graph matching networks for fewshot 3d action recognition. In *European Conference on Computer Vision (ECCV)*, pages 653–669, 2018.
- [19] F. Monti, D. Boscaini, J. Masci, E. Rodola, J. Svoboda, and M. M. Bronstein. Geometric deep learning on graphs and manifolds using mixture model cnns. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 5423–5434, 2017.
- [20] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng. Reading digits in natural images with unsupervised feature learning. In NIPS workshop on deep learning and unsupervised feature learning, 2011.
- [21] F. Nie, G. Cai, and X. Li. Multi-view clustering and semi-supervised classification with adaptive neighbours. In AAAI, pages 2408–2414, 2017.
- [22] F. Nie, J. Li, X. Li, et al. Parameter-free auto-weighted multiple graph learning: A framework for multiview clustering and semi-supervised classification. In *IJCAI*, pages 1881–1887, 2016.
- [23] B. Perozzi, R. Al-Rfou, and S. Skiena. Deepwalk: Online learning of social representations. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 701–710, 2014.
- [24] W. L. H. P. L. Y. B. Petar Veličković, William Fedus and R. D. Hjelm. Deep graph infomax. In ICLR, 2019.
- [25] M. Schlichtkrull, T. N. Kipf, P. Bloem, R. van den Berg, I. Titov, and M. Welling. Modeling relational data with graph convolutional networks. In *European Semantic Web Conference*, pages 593–607, 2018.
- [26] P. Sen, G. Namata, M. Bilgic, L. Getoor, B. Galligher, and T. Eliassi-Rad. Collective classification in network data. AI magazine, 29(3):93, 2008.
- [27] M. Simonovsky and N. Komodakis. Dynamic edge-conditioned filters in convolutional neural networks on graphs. In 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 29–38, 2017.
- [28] B. J. J. S. S.-C. Z. Siyuan Qi, Wenguan Wang. Learning human-object interactions by graph parsing neural networks. In *European Conference on Computer Vision (ECCV)*, pages 401–417, 2018.
- [29] G. Te, W. Hu, Z. Guo, and A. Zheng. Rgcnn: Regularized graph cnn for point cloud segmentation. arXiv preprint arXiv:1806.02952, 2018.
- [30] P. Velickovic, G. Cucurull, A. Casanova, A. Romero, P. Lio, and Y. Bengio. Graph attention networks. arXiv preprint arXiv:1710.10903, 2017.

- [31] J. Winn and N. Jojic. Locus: Learning object classes with unsupervised segmentation. In IEEE International Conference on Computer Vision, pages 756–763, 2005.
 [32] T. Xia, D. Tao, T. Mei, and Y. Zhang. Multiview spectral embedding. *IEEE Transactions on Systems*,
- Man, and Cybernetics, Part B (Cybernetics), 40(6):1438–1446, 2010.

 [33] Z. Xinyi and L. Chen. Capsule graph neural network. In *ICLR*, 2019.

 [34] X. Zhu, Z. Ghahramani, and J. D. Lafferty. Semi-supervised learning using gaussian fields and harmonic
- functions. In Proceedings of the 20th International conference on Machine learning (ICML-03), pages
- [35] C. Zhuang and Q. Ma. Dual graph convolutional networks for graph-based semi-supervised classification. In World Wide Web Conference on World Wide Web, pages 499–508, 2018.