

Graph Neural Networks - Notes

Nihal V. Nayak

Update: September 2020

Introduction

Graph Neural Networks (GNN) is a type of neural network which learns the structure of a graph. Learning graph structure allows us to represent the nodes of the graph in the euclidean space which can be useful for several downstream machine learning tasks. Recent work on GNN has shown impressive performance on link prediction, node classification and semi-supervised classification tasks. Since there is a growing interest in the machine learning community to learn more about these techniques, this document provides an introduction to GNN¹.

The document is organized as follows: First, we introduce the basic concepts of graphs and networks. Second, we describe the major steps used in GNNs to compute node embeddings. Next, we describe three GNN techniques that are commonly mentioned in the existing literature. Finally, we include a limited survey of other notable works in the field.

1 Basics

In this section, we will review the basics of networks that will be useful to understand the forthcoming sections.

Graphs consist of a set of nodes V and edges E . They are normally represented with an unweighted adjacency matrix A , which consists of 1s and 0s.

For example, take a look at the graph in figure 1. The graph has 3 nodes and 2 edges. The adjacency matrix is as follows:

$$\begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} \quad (1)$$

The degree matrix is a diagonal matrix that contains information about the degree of each node. The Degree of a node indicates the number of terminating edges for a node (self-loop counts as 2). For example, the degree matrix for fig. 1 is as follows:

¹This document assumes that you are familiar with the basics of deep learning

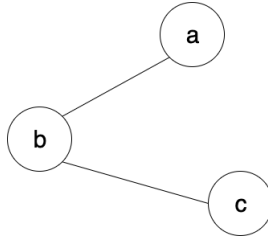


Figure 1: Graph with 3 nodes and 2 undirected edges

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad (2)$$

In GNN, each node v is associated with a feature vector $x_v \in \mathcal{R}^d$. Typically, the feature vector is randomly initialized and trained on a downstream task.

Although we have not introduced graph neural network, it would be useful to know how to compute the neighbors for a node. For example, from Figure 1 the neighbours of the node b are:

$$\mathcal{N}(b) = \{a, c\} \quad (3)$$

where $\mathcal{N}(\cdot)$ is the neighbour function. Likewise, the neighbor of the node a is:

$$\mathcal{N}(a) = \{b\} \quad (4)$$

As you can see, each node can have a varying number of neighbors. Therefore, we need a neural network that can deal with the varying number of neighbors.

2 Learning on Graphs

Graph neural network (GNN) is a family of algorithms that learns the structure of the graph in the euclidean space (Hamilton et al., 2017b).

A basic GNN consists of two components:

- **Aggregate:** For a given node, the Aggregate step applies a permutation invariant function to its neighbors to generate the aggregated node feature
- **Combine:** For a given node, Combine step passes the aggregated node feature to a learnable layer to generate node embedding for the GNN layer

The GNN layers are stacked together to increase the receptive field of clustering. For example, in Fig. 2 we have 2 GNN layers. This means that for every node, we aggregate over the 2-hop neighborhood.

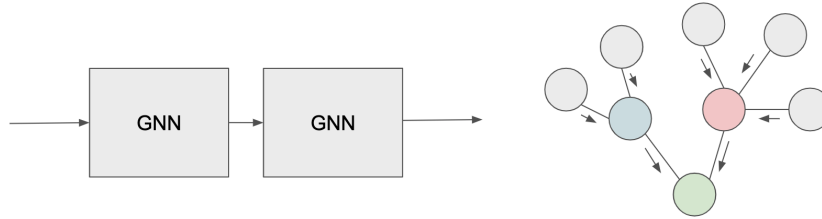


Figure 2: *Left*: 2 GNN layers stacked together and *Right*: shows the how the node feature is generated for a node (green) in a 2 layer GNN. The node (green) learns the node embedding by aggregating over 2-hop neighbourhood.

Aggregate and Combine components are also called message passing phase and readout phase in the message passing neural network framework (Gilmer et al., 2017).

2.1 Formal Definition

Consider the Graph $G = (V, E)$, where V is the set of node with features X_v . There are two main steps - Aggregate and Combine at each l -th layer of the GNN:

$$a_v^{(l)} = \text{Aggregate}^{(l)} \left(\left\{ h_u^{(l-1)} \forall u \in \mathcal{N}(v) \right\} \right) \quad (5)$$

where $a_v^{(l)}$ is the aggregated node feature of the neighbourhood, $h_u^{(l-1)}$ is the node feature in neighbourhood $\mathcal{N}(\cdot)$ of node v .

$$h_v^{(l)} = \text{Combine}^{(l)} \left(h_v^{(l-1)}, a_v^{(l)} \right) \quad (6)$$

where $h_v^{(l)}$ is the node representation at the l -th iteration. $h_v^{(0)} = x_v$ where x_v is the initial feature vector for the node.

3 Graph Convolutional Networks (GCN)

GCN (Kipf and Welling, 2016) is a graph neural network technique that makes use of the symmetrically normalized graph Laplacian to compute the node embeddings. Each GCN block receives node features from the $(l-1)$ th GCN block, i.e. the node features from the previous layer is passed to the next GCN layer, which allows for greater Laplacian smoothing (Li et al., 2018).

3.1 Aggregate

Since we don't always have access to the entire graph structure to compute the symmetrically normalized graph Laplacian, GCN has been interpreted as a mean aggregator (Xu et al., 2018a). Therefore,

$$a_v^{(l)} = \text{Mean} \left(\left\{ h_u^{(l-1)}, u \in \mathcal{N}(v) \cup \{v\} \right\} \right) \quad (7)$$

For each node v , the algorithm takes the mean of the neighborhood node features along with its feature (self-loop) to compute the aggregated node feature $a_v^{(l)}$ for the l th layer. This formulation makes GCN an inductive graph neural network.

3.2 Combine

The aggregated node feature $a_v^{(l)}$ is multiplied with a learnable weight matrix followed by a non-linear activation (ReLU) to get the node feature $h_v^{(l)}$ for the l th layer,

$$h_v^{(l)} = \text{ReLU} \left(W^{(l)} a_v^{(l)} \right) \quad (8)$$

where $W^{(l)} \in R^{d_{(l-1)} \times d_{(l)}}$ is the learnable weight matrix.

4 GraphSage

GraphSage (Hamilton et al., 2017a) is a spatial graph neural network algorithm that learns node embeddings by sampling and aggregating neighbors from multiple hops. GraphSage, unlike GCN, is an inductive learning algorithm which means that it does not require the entire graph structure during learning.

In the original GraphSage paper, the authors describe multiple aggregators - mean-pooling, max-pooling and LSTM aggregator. In this section, we'll be describing the LSTM aggregator and the combine from the original paper. For simplicity, we do not sample the neighborhood nodes as well.

4.1 Aggregate

GraphSage makes use of a LSTM aggregator to compute an aggregated node feature. Since LSTM is not permutation invariant, the authors shuffle the neighbourhood nodes to compute the aggregated vector.

$$a_v^{(l)} = \text{LSTM} \left(\left\{ h_u^{(l-1)}, u \in \mathcal{N}(v) \right\} \right) \quad (9)$$

where $a_v^{(l)}$ is the aggregated node feature for the l th layer. $W_{\text{pool}}^{(l-1)}$ and b are the learnable parameters.

4.2 Combine

The combine function concatenates the $(l-1)$ th layer node feature $h_v^{(l-1)}$ with the aggregated feature vector $a_v^{(l)}$ and pass the vector through a fully connected neural network:

$$h_v^{(l)} = \text{ReLU} \left(W^{(l)} [h_v^{(l-1)} || a_v^{(l)}] \right) \quad (10)$$

where $W^{(l)} \in R^{d_{(l-1)} \times d_{(l)}}$ is a learnable weight matrix for l th layer.

5 Graph Attention Network (GAT)

Graph Attention Network (Veličković et al., 2018) is a spatial graph neural network technique that uses Self Attention to aggregate the neighborhood node features. Self Attention is equivalent to computing a weighted mean of the neighbourhood node features.

The original paper uses K *masked self attention* modules to aggregate node features. For simplicity, we consider only one self attention block.

5.1 Aggregate

GAT uses an Attention module to compute the weighted mean of the neighborhood node features as follows:

$$\alpha_u^{(l)} = \text{Attention} \left(\left\{ (Wh_u^{(l-1)} || Wh_v^{(l-1)}), u \in \mathcal{N}(v) \cup \{v\} \right\} \right) \quad (11)$$

where $\text{Attention}(\cdot)$ is single fully connected layer with LeakyReLU, $W \in R^{d_{(l-1)} \times d_{(l)}}$ and $\alpha_u^{(l)}$ is the coefficient for neighbour node u . The authors use weighted linear combination of neighbours (Bahdanau et al., 2014) but the aggregation technique is agnostic to the choice attention mechanism.

5.2 Combine

The aggregated node feature $a_v^{(l)}$ is multiplied with a learnable weight matrix followed by a non-linear activation to get the node feature $h_v^{(l)}$ for the l th layer,

$$h_v^{(l)} = \text{ReLU} \left(\text{Sum} \left\{ \alpha_u^{(l)} Wh_u^{(l-1)}, u \in \mathcal{N}(v) \cup \{v\} \right\} \right) \quad (12)$$

6 Related Works

Apart from GCN, Graphsage, and GAT, there are few more works that frequently appear in the literature. SGN (Wu et al., 2019) simplifies GCN by removing non-linearity between the GCN layers. Furthermore, to achieve the same “receptive field” of having K GCN layers, they take the graph laplacian to the power of K and show equivalent performance. Xu et al. (2018b) introduced jumping knowledge networks that aggregate node features generated from all the previous layers to compute the final vector using mean or max-pooling. This allows the model to learn better structure-aware representations of the graph. R-GCN (Schlichtkrull et al., 2018) and Directed-GCN (Marcheggiani and Titov, 2017) extend the GCN to relational graphs.

References

- Bahdanau, D., Cho, K., and Bengio, Y. (2014). Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*.
- Gilmer, J., Schoenholz, S. S., Riley, P. F., Vinyals, O., and Dahl, G. E. (2017). Neural message passing for quantum chemistry. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1263–1272. JMLR. org.
- Hamilton, W., Ying, Z., and Leskovec, J. (2017a). Inductive representation learning on large graphs. In *Advances in Neural Information Processing Systems*, pages 1024–1034.
- Hamilton, W. L., Ying, R., and Leskovec, J. (2017b). Representation learning on graphs: Methods and applications. *IEEE Data Eng. Bull.*, 40:52–74.
- Kipf, T. N. and Welling, M. (2016). Semi-supervised classification with graph convolutional networks.
- Li, Q., Han, Z., and Wu, X.-M. (2018). Deeper insights into graph convolutional networks for semi-supervised learning. In *Thirty-Second AAAI Conference on Artificial Intelligence*.
- Marcheggiani, D. and Titov, I. (2017). Encoding sentences with graph convolutional networks for semantic role labeling. *arXiv preprint arXiv:1703.04826*.
- Schlichtkrull, M., Kipf, T. N., Bloem, P., Van Den Berg, R., Titov, I., and Welling, M. (2018). Modeling relational data with graph convolutional networks. In *European Semantic Web Conference*, pages 593–607. Springer.
- Veličković, P., Cucurull, G., Casanova, A., Romero, A., Liò, P., and Bengio, Y. (2018). Graph Attention Networks. *International Conference on Learning Representations*.
- Wu, F., Souza, A., Zhang, T., Fifty, C., Yu, T., and Weinberger, K. (2019). Simplifying graph convolutional networks. In *International Conference on Machine Learning*, pages 6861–6871.
- Xu, K., Hu, W., Leskovec, J., and Jegelka, S. (2018a). How powerful are graph neural networks?
- 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*.

Appendix

A Graph Laplacian

The adjacency matrix and degree matrix can be combined to get the graph Laplacian.

$$L = D - A \quad (13)$$

where L is the graph laplacian, D is the degree matrix and A is the adjacency matrix.

The symmetric normalized graph laplacian can be formally defined as:

$$L_{\text{sym}} = D^{-\frac{1}{2}} L D^{-\frac{1}{2}} = I - D^{-\frac{1}{2}} A D^{-\frac{1}{2}} \quad (14)$$

where I is the Identity matrix.

The elements of L_{sym} are given by:

$$L_{ij} = \begin{cases} 1 & \text{if } i == j \text{ and } deg(v_i) \neq 0 \\ -\frac{1}{\sqrt{deg(v_i)deg(v_j)}} & \text{if } i \neq j \\ 0 & \text{otherwise} \end{cases} \quad (15)$$

B Original GCN

In section 3, we interpreted GCN (Kipf and Welling, 2016) as a mean aggregator. However, there is a subtle difference between our interpretation and the original GCN.

Consider the adjacency matrix A and its degree matrix D . Let $H^{(l)}$ be the node features for the l th layer and $H^{(0)} = X$. The propagation rule for each GCN layer is as follows:

$$f(H^{(l)}, \tilde{A}) = \sigma \left(\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} H^{(l)} W^{(l)} \right) \quad (16)$$

with $\tilde{A} = I + A$ where I is the identity matrix and degree matrix \tilde{D} of \tilde{A} .