
Optimizing Graph Neural Networks for Enhanced Community Detection in Social Networks

Parsa Kamalipour

Department of Computer Science

Concordia University

parsa.kamalipour@mail.concordia.ca

Abdulmoumen Al-Atrash

Department of Computer Science

Concordia University

a_alatr@live.concordia.ca

Aniket Roy

Department of Computer Science

Concordia University

aniket.roy@mail.concordia.ca

Abstract

Community detection, a key problem in social network analysis, identifies clusters in graphs where nodes are more interconnected than with the rest of the network. This has applications in fields like social network analysis and biology, aiding in understanding complex systems. The rapid growth of social networks has made analyzing these intricate structures increasingly critical. Traditional methods, such as modularity optimization and spectral clustering, struggle with large networks and often miss finer community structures. Optimization-based approaches like Learning-Based Genetic Algorithms (LGA) and Quantum-Inspired Optimization show high accuracy but lack generalization across diverse networks. Modularity-based methods, while efficient, rely heavily on modularity maximization, limiting their effectiveness in distinguishing communities in networks of varying sizes. Deep learning, particularly Graph Neural Networks (GNNs), has proven effective in addressing these challenges, offering scalability, accuracy, and adaptability to diverse networks. This project leverages GNNs to overcome these limitations by learning high-dimensional feature representations of nodes and communities, improving the accuracy and scalability of community detection. It aligns with course objectives by applying theoretical knowledge to practical problems in network analysis and machine learning, addressing gaps in traditional methods with enhancements like graph sampling, optimized graph convolutions, and graph partitioning.

1 Introduction

1.1 Introduction to Community Detection

Community detection is all about identifying clusters or groups within a network - which are called communities - where nodes are more densely connected to each other than to the rest of the network. This problem is very crucial in various domains, including social network analysis, biology, and information retrieval, as it helps uncover the underlying structure and functional organization of complex systems. For instance, community detections in biology aid in identifying interconnected gene networks, while in social networks, they reveal social group dynamics and interaction patterns. Traditional methods, such as modularity optimization and spectral clustering, have been widely used for this purpose. However, these approaches often face challenges in handling large-scale

networks and capturing intricate community structures, particularly when networks exhibit complex relationships or irregular structural patterns. As social networks continue to grow and become more deeply integrated into daily life, it is becoming vital to analyze these complex structures efficiently and on a large scale. This development demands advanced solutions that maintain performance while achieving scalability to keep up with the diverse and dynamic nature of real-world networks.

1.2 Machine Learning Context

Recent advancements in machine learning, particularly deep learning, have introduced numerous methodologies for community detection. Graph Neural Networks (GNNs) have proven to be powerful tools for learning complex patterns in graph-structured data. For example, the Contrastive Deep Non-negative Matrix Factorization (CDNMF) model combines contrastive learning with deep nonnegative matrix factorization, capturing both network topology and node attributes to enhance community detection [1]. Similarly, integrating community detection algorithms with GNNs has improved link prediction in scientific literature networks, showcasing the synergy of these approaches [2]. This underscores the potential of GNNs to overcome the limitations of modularity-based and optimization-based methods, offering greater accuracy and scalability for complex networks.

1.3 Usage of Machine Learning in Community Detection

In this project, we aim to use machine learning techniques to address the challenges of community detection in large-scale networks, to be more precise in Social Networks. By using models such as GNNs and incorporating methods like optimization on previous ideas, we seek to develop an algorithm that can effectively identify communities with high accuracy and scalability. This approach aligns with the course objectives by applying machine learning methodologies to solve complex problems; thereby, by doing this research, we both learn a new Machine Learning methodology and also enhance our understanding of community structures in various real-world applications.

1.4 Challenges in Community Detection

Traditional approaches to community detection face notable limitations, particularly in modern large-scale applications. Scalability is a key challenge, as real-world networks often contain millions or billions of nodes and edges, requiring methods capable of efficient processing [3]. Another issue is the noisy and irregular nature of network data [4]. False positives (incorrect links) and false negatives (missing links) can distort community structures, reducing accuracy for less sophisticated methods. Additionally, heterogeneous networks, where nodes and edges carry varying data types, pose challenges for algorithms reliant on structural relationships to capture context [5]. Community detection is inherently indeterminate, often lacking a single correct solution, demanding robust methods that generalize across diverse network structures and sizes. These challenges form the foundation for the discussions in this paper.

2 Related Work

2.1 Overview

This section presents a literature review to examine pre-existing knowledge in the area of the NP-Hard problem of 'Community Detection' in social networks. We examine relevant existing work to ensure that our proposed approach is justified based on the problem context and that it differs from already proposed solutions to the NP-Hard problem. In this section, we identify and discuss three main categories of algorithms used to tackle the problem of community detection, namely optimization-based methods, modularity-based approaches, and deep learning techniques.

2.2 Optimization-Based Methods

Optimization is a popular approach for tackling complex machine learning problems like community detection. The Learning-Based Genetic Algorithm (LGA)[6] [Identifying Communities in Complex Networks Using Learning-Based Genetic Algorithm] addresses the NP-hard nature of community detection by combining Learning Automata with genetic operations, overcoming issues like premature

convergence and local optima. LGA achieves a 26.47% accuracy improvement on real-world networks and a 48.32% improvement on synthetic networks. Akbar et al.[7] introduce a quantum mechanics-inspired optimization method for detecting patterns in ecological communities, simulating quantum principles to identify biodiversity changes caused by environmental factors like land use and climate change. MARLCD [8], a multi-agent reinforcement learning algorithm, uses agents to independently explore network communities, outperforming methods like GA-Net and Meme-Net on real and synthetic datasets. However, these approaches lack generalizability across networks of varying domains and complexities. Their assumptions, while effective in specific contexts, limit flexibility and can lead to excessive computational costs when applied to networks with diverse structural patterns.

2.3 Modularity-Based Methods

Modularity-based algorithms are simple and interpretable techniques widely used for community detection by maximizing modularity. Modularity measures the density of connections within communities compared to their connections with the rest of the network. Chen et al.[9] propose a hierarchical clustering algorithm to optimize Max-Min Modularity, arguing that traditional statistical inference assumes data independence and relation-based methods fail to capture domain-specific features, leading to lower performance. The Hybrid Genetic Tabu (HGT)[10] method combines Genetic Algorithms (GA) and Tabu Search (TS) to maximize modularity, outperforming methods like Louvain and Label Propagation, especially in medium-sized networks. However, modularity-based approaches rely heavily on modularity maximization, which often struggles to differentiate between community sizes and connectivity, resulting in suboptimal detection accuracy.

2.4 Deep Learning Methods

Deep learning is often the go-to choice for handling complex data processing, feature extraction, and pattern recognition. Liu et al.[11] explore various deep learning techniques for community detection, including deep neural networks, graph neural networks (GNNs), and deep graph embeddings. The paper highlights GNNs' ability to encode higher-dimensional feature representations of nodes and communities, addressing limitations of traditional methods that rely heavily on adjacency and attribute matrices, often losing complex structural relationships. The deep transitive encoder[12] transforms the adjacency matrix to capture indirect node relationships missed by traditional methods. Using unsupervised transfer learning, it extracts low-dimensional features, achieving improved accuracy. However, its reliance on adjacency matrix transformation makes it computationally demanding for large-scale networks. Nooribakhsh et al. [13] provide a systematic overview of machine learning trends in community detection, emphasizing the rise of deep learning in tackling large, complex networks. They highlight deep learning's consistent outperformance of simpler methods like game-theoretic and clustering algorithms, as well as its superiority over density-based approaches, which struggle to capture structural and feature-based characteristics of nodes.

3 Proposed Approach

Scalability is a key challenge in community detection, especially for large-scale graphs. Traditional methods often struggle with the computational demands of modern networks, leading to high memory usage and prolonged training times. This study evaluates the scalability of Graph Neural Networks (GNNs) through three methodologies: full-batch training, neighbor sampling, and graph partitioning. These strategies aim to improve computational efficiency while maintaining model accuracy.

The main objective is to analyze and compare methods for addressing scalability issues in large graphs during community detection. Training or evaluating the entire graph at once requires substantial memory, particularly for large graphs. To mitigate this, we compare multiple approaches to determine which offers the best balance between accuracy and training time.

We implement models using GCNConv and GraphSAGEConv architectures, detailed in the accompanying code. The models are evaluated using the three training approaches discussed later in this section.

3.1 GNN Architectures

This study employs two well-established GNN architectures, **GCNConv** and **GraphSAGEConv**, as the foundational models:

3.1.1 GCNConv

The GCNConv class implements the graph convolutional operator from the paper "*Semi-supervised Classification with Graph Convolutional Networks*" by Kipf and Welling (2016) [14]. This operator performs message passing by aggregating feature information from a node's neighbors, facilitating the learning of node representations that capture both local graph structure and node features.

Mathematical Formulation: The layer computes the output features \mathbf{X}' as:

$$\mathbf{X}' = \hat{\mathbf{D}}^{-1/2} \hat{\mathbf{A}} \hat{\mathbf{D}}^{-1/2} \mathbf{X} \Theta, \quad (1)$$

where:

- $\hat{\mathbf{A}} = \mathbf{A} + \mathbf{I}$ is the adjacency matrix with added self-loops.
- $\hat{\mathbf{D}}$ is the diagonal degree matrix of $\hat{\mathbf{A}}$.
- \mathbf{X} denotes the input feature matrix.
- Θ represents the learnable weight matrix.

In a node-wise formulation, the output feature \mathbf{x}'_i for node i is:

$$\mathbf{x}'_i = \Theta^\top \sum_{j \in \mathcal{N}(i) \cup \{i\}} \frac{e_{j,i}}{\sqrt{\hat{d}_j \hat{d}_i}} \mathbf{x}_j, \quad (2)$$

where:

- $\mathcal{N}(i)$ denotes the set of neighbors of node i .
- $e_{j,i}$ is the edge weight from node j to node i (defaulting to 1.0 if not specified).
- $\hat{d}_i = 1 + \sum_{j \in \mathcal{N}(i)} e_{j,i}$ is the degree of node i in $\hat{\mathbf{A}}$.

3.1.2 GraphSAGEConv

The GraphSAGEConv class implements the GraphSAGE operator from the paper "*Inductive Representation Learning on Large Graphs*" by Hamilton et al. (2017) [15]. This operator enables inductive learning by aggregating feature information from a node's local neighborhood, allowing the model to generalize to unseen nodes or graphs.

Mathematical Formulation: The layer computes the output features \mathbf{x}'_i for node i as:

$$\mathbf{x}'_i = \mathbf{W}_1 \mathbf{x}_i + \mathbf{W}_2 \cdot \text{AGGREGATE}_{j \in \mathcal{N}(i)} \mathbf{x}_j, \quad (3)$$

where:

- \mathbf{W}_1 and \mathbf{W}_2 are learnable weight matrices.
- $\mathcal{N}(i)$ denotes the set of neighbors of node i .
- AGGREGATE is a function such as mean, LSTM, or max pooling that combines the features of neighboring nodes.

If the project parameter is set to True, each neighbor's feature \mathbf{x}_j is first transformed via a non-linear function:

$$\mathbf{x}_j \leftarrow \sigma(\mathbf{W}_3 \mathbf{x}_j + \mathbf{b}), \quad (4)$$

where:

- \mathbf{W}_3 is a learnable weight matrix.
- \mathbf{b} is a bias term.
- σ is an activation function (e.g., ReLU).

3.2 Scalability Techniques

3.2.1 Full-Batch Training

Full-batch training involves processing the entire graph simultaneously during each training epoch. While this method captures the global structure and relationships within the graph, it is computationally expensive for large graphs. Memory constraints often limit its applicability, as demonstrated by the prohibitive resource requirements for the Reddit dataset.

3.2.2 Neighbor Sampling

Neighbor sampling addresses the limitations of full-batch training by selectively sampling a fixed number of neighbors for each node during training. This approach significantly reduces memory usage while maintaining high accuracy. The methodology is based on the GraphSAGE architecture, which aggregates information from sampled neighbors to compute node embeddings.

3.2.3 Graph Partitioning

Graph partitioning divides the graph into smaller, non-overlapping subgraphs, which are processed independently during training. This method, inspired by the approach outlined in Chiang et al. [16], facilitates parallel processing and reduces the memory footprint. By isolating subgraphs, this technique effectively balances computational efficiency and model performance.

4 Results, Training, and Evaluation

In this section, we present the results of our experiments, including the training process, evaluation metrics, and analysis of the performance of the proposed approaches.

4.1 Datasets

We conducted experiments using three datasets of varying sizes and complexity to evaluate the scalability and performance of the proposed methods:

- **Synthetic Graphs:** Generated using the Stochastic Block Model (SBM), these graphs provide controlled environments to simulate community structures and test the model under predefined conditions.
- **CORA Dataset:** A citation network widely used in graph learning tasks, with moderate size and complexity.
- **Reddit Dataset:** A large-scale dataset with approximately 100x the nodes of the CORA dataset, highlighting scalability challenges.

4.2 Training Approaches

The training process was carried out using three distinct methodologies:

1. **Full-Batch Training:** The entire graph is processed simultaneously. This method faced memory constraints on large datasets like Reddit.
2. **Neighbor Sampling:** A subset of neighbors for each node is sampled during training, significantly reducing memory usage and enabling training on larger graphs.
3. **Graph Partitioning:** The graph is divided into smaller subgraphs processed independently, facilitating parallel computation and scalability.

4.3 Evaluation Metrics

To assess the effectiveness of the models and approaches, we used the following metrics:

- **Accuracy:** The ability of the model to correctly classify nodes into communities.
- **Training Time:** The computational time required to complete the training process.

- **Memory Usage:** The memory consumption during training, indicating the scalability of each method.

4.4 Experimental Results

4.4.1 Synthetic Graphs

On synthetic graphs, all methods performed well in terms of accuracy. Neighbor Sampling and Graph Partitioning showed substantial reductions in memory usage and training time compared to Full-Batch Training.

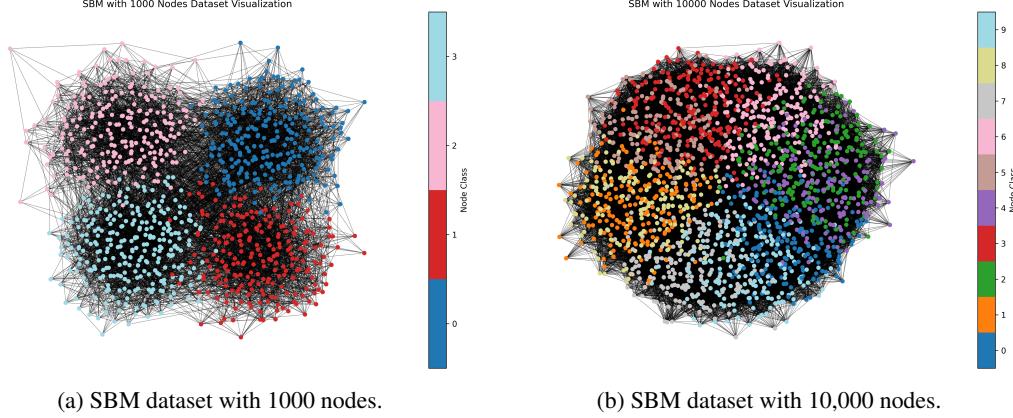


Figure 1: Visualization of SBM datasets with 1000 and 10,000 nodes. Node classes are color-coded to highlight community structures.

4.4.2 CORA Dataset

On the CORA dataset, Neighbor Sampling achieved comparable accuracy to Full-Batch Training but with reduced memory and time requirements. Graph Partitioning demonstrated similar efficiency gains while maintaining high accuracy.

4.4.3 Reddit Dataset

On the large-scale Reddit dataset, Full-Batch Training was infeasible due to excessive memory requirements. Both Neighbor Sampling and Graph Partitioning enabled training and achieved competitive accuracy. Among the two, Neighbor Sampling had slightly better training times, while Graph Partitioning excelled in memory optimization.

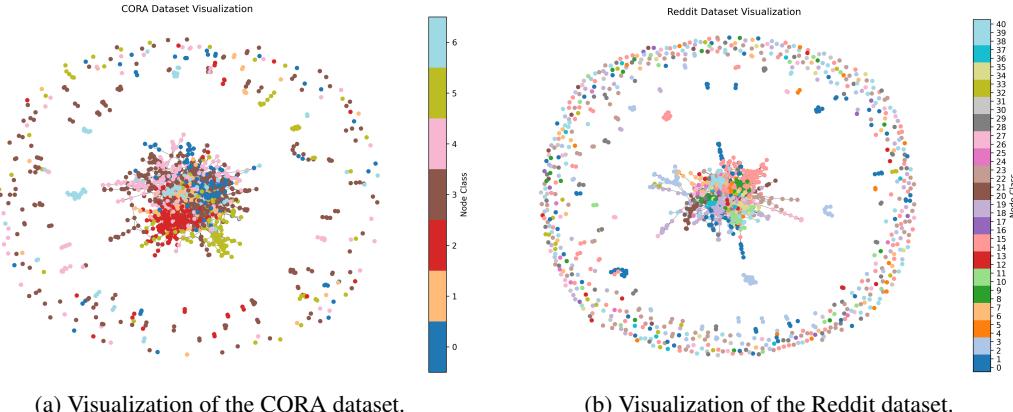


Figure 2: Visualizations of the CORA and Reddit datasets. Node classes are color-coded to represent different communities or research topics.

4.4.4 Performance Summary Tables

Table 1: Performance on SBM Dataset with 1,000 Nodes.

| Model | Method | Test Accuracy | Train Time (s) |
|--------------|--------------------|---------------|----------------|
| GCNet | Full Batch | 70% | 0.19 |
| | Neighbor Sampling | 67% | 3.32 |
| | Graph Partitioning | 50.5% | 0.58 |
| GraphSAGENet | Full Batch | 52.5% | 0.13 |
| | Neighbor Sampling | 85.5% | 3.17 |
| | Graph Partitioning | 46% | 0.65 |

Table 2: Performance on SBM Dataset with 10,000 Nodes.

| Model | Method | Test Accuracy | Train Time (s) |
|--------------|--------------------|---------------|----------------|
| GCNet | Full Batch | 90% | 12.66 |
| | Neighbor Sampling | 17.85% | 44.62 |
| | Graph Partitioning | 89.85% | 9.51 |
| GraphSAGENet | Full Batch | 20.7% | 10.23 |
| | Neighbor Sampling | 24.05% | 42.07 |
| | Graph Partitioning | 58% | 8.40 |

Table 3: Performance on CORA Dataset.

| Model | Method | Test Accuracy | Train Time (s) |
|--------------|--------------------|---------------|----------------|
| GCNet | Full Batch | 89.48% | 0.19 |
| | Neighbor Sampling | 88.01% | 4.86 |
| | Graph Partitioning | 87.08% | 0.86 |
| GraphSAGENet | Full Batch | 90.41% | 0.88 |
| | Neighbor Sampling | 88.56% | 8.93 |
| | Graph Partitioning | 88.75% | 1.58 |

Table 4: Performance on Reddit Dataset.

| Model | Method | Test Accuracy | Train Time (s) |
|--------------|--------------------|---------------------|----------------|
| GCNet | Full Batch | N/A (Out of Memory) | N/A |
| | Neighbor Sampling | 88.6% | 1892.33 |
| | Graph Partitioning | 88.12% | 1688.57 |
| GraphSAGENet | Full Batch | N/A (Out of Memory) | N/A |
| | Neighbor Sampling | 86.44% | 13108.7 |
| | Graph Partitioning | N/A (Out of Memory) | N/A |

5 Conclusion

In conclusion, this study demonstrates the effectiveness of scalable Graph Neural Network (GNN) training methods for addressing the challenges of community detection in large-scale graphs. Neighbor Sampling offers a balance between accuracy and computational efficiency, while Graph Partitioning provides superior memory optimization, emphasizing the importance of selecting the appropriate method based on dataset characteristics and resource constraints. By leveraging GNNs, this work overcomes the limitations of traditional modularity-based and optimization approaches, which often lack generalizability and scalability. The proposed methodologies enhance the practicality of community detection in large, complex networks, bridging the gap between theoretical advancements and real-world applications. These findings not only contribute to a deeper understanding of network structures but also highlight the potential of advanced deep learning techniques to tackle the NP-hard problem of community detection effectively.

References

- [1] Yuecheng Li, Jialong Chen, Chuan Chen, Lei Yang, and Zibin Zheng. Contrastive deep nonnegative matrix factorization for community detection, 2024.
- [2] Chunjiang Liu, Yikun Han, Haiyun Xu, Shihan Yang, Kaidi Wang, and Yongye Su. A community detection and graph neural network based link prediction approach for scientific literature, 2024.
- [3] Xinyu Que, Fabio Checconi, Fabrizio Petrini, and John A. Gunnels. Scalable community detection with the louvain algorithm. In *2015 IEEE International Parallel and Distributed Processing Symposium*, pages 28–37, 2015.
- [4] Yiye Ruan, David Fuhr, and Srinivasan Parthasarathy. Efficient community detection in large networks using content and links. In *Proceedings of the 22nd International Conference on World Wide Web, WWW ’13*, page 1089–1098, New York, NY, USA, 2013. Association for Computing Machinery.
- [5] Chuan Shi, Yitong Li, Jiawei Zhang, Yizhou Sun, and Philip S. Yu. A survey of heterogeneous information network analysis. *IEEE Transactions on Knowledge and Data Engineering*, 29(1):17–37, 2017.
- [6] Gholam Reza Abdi, Amir Hosein Refahi Sheikhani, Sohrab Kordrostami, Bagher Zarei, and Mohsen Falah Rad. Identifying communities in complex networks using learning-based genetic algorithm. *Ain Shams Engineering Journal*, page 103031, 2024.
- [7] S. Akbar and S. K. Saritha. Quantum inspired community detection for analysis of biodiversity change driven by land-use conversion and climate change. *Scientific Reports*, 11(1):14332, 2021.
- [8] Mir Mohammad Alipour and Mohsen and Abdolhosseinzadeh. A multiagent reinforcement learning algorithm to solve the community detection problem. *Signal and Data Processing*, 19(1), 2022.
- [9] Jiyang Chen, Osmar R. Zaïane, and Randy Goebel. *Detecting Communities in Social Networks using Max-Min Modularity*, pages 978–989.
- [10] Bouchéma Sara Cheikh Salmi and Zaoui Sara. An enhanced evolutionary approach for solving the community detection problem. *Journal of Information and Telecommunication*, 6(1):83–100, 2022.
- [11] Fanzhen Liu, Shan Xue, Jia Wu, Chuan Zhou, Wenbin Hu, Cecile Paris, Surya Nepal, Jian Yang, and Philip S. Yu. Deep learning for community detection: progress, challenges and opportunities. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI’20*, 2021.
- [12] Ying Xie, Xinmei Wang, Dan Jiang, and Rongbin Xu. High-performance community detection in social networks using a deep transitive autoencoder. *Information Sciences*, 493:75–90, 2019.
- [13] Mahsa Nooribakhsh, Marta Fernández-Diego, Fernando González-Ladrón-De-Guevara, and Mahdi Mollamotalebi. Community detection in social networks using machine learning: a systematic mapping study. *Knowledge and Information Systems*, 66(12):7205–7259, December 2024.
- [14] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks, 2017.
- [15] William L. Hamilton, Rex Ying, and Jure Leskovec. Inductive representation learning on large graphs, 2018.
- [16] Wei-Lin Chiang, Xuanqing Liu, Si Si, Yang Li, Samy Bengio, and Cho-Jui Hsieh. Cluster-gcn: An efficient algorithm for training deep and large graph convolutional networks. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery and Mining, KDD ’19*. ACM, July 2019.