# Community Detection in Social Network Analysis

*by*

## Aman Kumar Sharma

## (22MCB0021)

## GitHub Link:

# INDEX

# 1. ABSTRACT:

The web's expansion and emergence of many social networking sites (SNS) have empowered users to easily interconnect on a shared platform. A social network can be represented by a graph consisting of a set of nodes and edges connecting these nodes. The nodes represent the individuals/entities, and the edges correspond to the interactions among them. The tendency of people with similar tastes, choices, and preferences to get associated with a social network leads to the formation of virtual clusters or communities. Detection of these communities can be beneficial for numerous applications such as finding a common research area in collaboration networks, finding a set of like-minded users for marketing and recommendations, and finding protein interaction networks in biological networks. Many community-detection algorithms have been proposed and applied in this project. This paper presents a survey of the existing algorithms and approaches for the detection of communities in social networks. We also discuss some of the applications of community detection.

# 2. INTRODUCTION:

Community detection algorithms play a pivotal role in social network analysis, aiding in the exploration and understanding of complex social structures. Social networks are composed of individuals or entities connected by various relationships, such as friendships, collaborations, or interactions. By identifying communities within these networks, community detection algorithms unveil groups of nodes that exhibit strong interconnections and share similar characteristics or behaviors. The study of communities within social networks has garnered significant attention due to its applications in various fields. Sociologists are interested in understanding social structures, identifying influential individuals, and studying information diffusion. Biologists utilize community detection to analyze protein-protein interaction networks and identify functional modules within biological systems. In marketing, community detection helps target specific customer segments and design tailored advertising campaigns. Additionally, in online social networks, detecting communities facilitates recommendation systems, content personalization, and understanding of user behavior. Community detection algorithms employ diverse methodologies to uncover communities within networks. Some algorithms utilize graph-based approaches, leveraging network properties such as node degree,

network density, and clustering coefficients. Others employ optimization techniques to maximize a specific quality function, such as modularity or conductance, to partition the network into cohesive communities. Furthermore, algorithms can be hierarchical, dividing communities into nested sub-communities, or non-hierarchical, directly assigning nodes to communities. Despite their wide application, community detection algorithms face challenges. The resolution limit problem arises when algorithms fail to detect small communities within large-scale networks. Overlapping communities, where nodes can belong to multiple communities simultaneously, pose additional complexity. Furthermore, the choice of algorithm and parameter settings can influence the accuracy and effectiveness of community detection results. In recent years, advancements in network analysis techniques, computational power, and data availability have fueled research in community detection algorithms. Researchers are continuously developing novel algorithms that address scalability issues, detect overlapping communities, and consider dynamic networks' evolving nature. Furthermore, interdisciplinary collaborations are prevalent, as researchers from various domains contribute to advancing community detection algorithms and applying them to real-world scenarios.

This study provides an overview of community detection algorithms in social network analysis, exploring their methodologies, applications, challenges, and ongoing research directions. By uncovering the underlying structure of communities, these algorithms enhance our understanding of complex social systems and contribute to a wide range of disciplines.

# 3. DATASET OVERVIEW:

**Dolphin Social Network Dataset:**

https://www.kaggle.com/datasets/mashazhil/social-network-of-dolphins-in-new-zealand

The Dolphin Social Network dataset is a well-known dataset used in social network analysis. It consists of a network of interactions between bottlenose dolphins (Tursiops truncatus) living in Doubtful Sound, New Zealand. The dataset captures the social relationships and interactions among the dolphins, providing valuable insights into their social behavior and dynamics.

Here is an overview of the Dolphin Social Network dataset:

1. Origin: The dataset was collected through field observations and recordings of dolphin interactions in Doubtful Sound, New Zealand. Researchers spent time observing and documenting the dolphins' behaviors, social interactions, and associations.

2. Nodes: Each node in the network represents an individual dolphin. The dolphins are identified by unique IDs or names, allowing researchers to track their behavior over time.

3. Edges: The edges in the network represent social interactions between dolphins. An edge between two dolphins indicates that they have been observed engaging in social behaviors such as swimming together, vocalizing, playing, or exhibiting affiliative behaviors.

4. Weighted or Unweighted: The dataset may include information about the strength or frequency of interactions between dolphins, making it a weighted network. Alternatively, it could be an unweighted network where the presence of an edge signifies the existence of any social interaction, without considering its intensity.

5. Network Structure: The Dolphin Social Network dataset exhibits a typical social network structure with patterns of connectivity and clustering. It is likely to have core individuals with many connections, as well as peripheral individuals with fewer connections. The presence of communities or subgroups within the dolphin population may also be observed.

6. Temporal Dynamics: Depending on the data collection methodology, the dataset may capture the temporal aspect of dolphin interactions. This allows researchers to analyze the evolution of social relationships over time, detect changes in network structure, and investigate the influence of external factors on the social network.

7. Research Applications: The Dolphin Social Network dataset has been utilized in various research studies in fields such as animal behavior, social network analysis, ecology, and marine biology. Researchers have examined topics like social structure, the centrality of individuals, communication patterns, social learning, and the impact of environmental factors on social networks.

The Dolphin Social Network dataset offers a valuable resource for studying social interactions and relationships among dolphins. It provides researchers with a real-

world network to explore and analyze using social network analysis techniques, shedding light on the complex social dynamics of these intelligent marine mammals.

# 4. ALGORITHMS OVERVIEW AND INTRODUCTION:

Here I am using CPM (Clique Percolation Method), Louvain Modularity Algorithm, and Girvan-Newman Algorithm for community detection in social network analysis. Each algorithm has its unique approach and characteristics for identifying communities within networks. Here is an overview of these algorithms:

## 1. Clique Percolation Method (CPM):

The CPM algorithm, also known as the Clique Percolation Method, is a community detection algorithm used to identify overlapping communities in a network. It was proposed by Palla, Derényi, Farkas, and Vicsek in 2005. The algorithm is based on the concept of k-cliques, which are complete subgraphs of size k in the network.

Here's a step-by-step overview of the CPM algorithm:

1. Define the parameter k: The algorithm requires specifying the size of the cliques, denoted by k. A clique of size k is a fully connected subgraph with k nodes.

2. Find all k-cliques: Generate a list of all k-cliques in the network. This can be done using algorithms such as the Bron-Kerbosch algorithm.

3. Build a clique graph: Create a clique graph where each node represents a unique k-clique, and there is an edge between two nodes if the corresponding cliques share k-1 nodes.

4. Find the overlapping communities: Identify the communities in the clique graph using a percolation process. Initially, each node represents a separate community. Merge communities if they share k-1 nodes or more.

5. Post-process the communities: Remove any duplicate or overlapping communities that may have been formed during the percolation process. This step ensures that each node belongs to only one community.

The CPM algorithm allows for the identification of overlapping communities because each node can be a member of multiple communities. The algorithm's effectiveness depends on choosing an appropriate value for k, which determines the size of the cliques. Larger values of k lead to more specific communities, while smaller values may result in more overlapping communities. It's worth noting that the CPM algorithm is computationally expensive, particularly when dealing with large networks. Various optimizations and heuristics have been proposed to improve its efficiency and scalability, such as reducing the number of k-cliques considered or using parallelization techniques.

Overall, the CPM algorithm provides a framework for detecting overlapping communities in networks based on the concept of k-cliques and percolation.

## 2. Louvain Modularity Algorithm:

he Louvain Modularity Algorithm is a popular community detection algorithm known for its efficiency and ability to identify communities with high modularity in large-scale networks. It was developed by Vincent D. Blondel, Jean-Loup Guillaume, Renaud Lambiotte, and Etienne Lefebvre in 2008.

Here is an overview of the Louvain Modularity Algorithm:

1. Modularity Metric:

The algorithm optimizes a quality function called modularity to identify communities within a network. Modularity measures the degree to which the network is divided into communities based on the density of edges within communities compared to a null model. Higher modularity values indicate a stronger division of the network into communities.

2. Algorithm Steps:

Initially, each node is considered as a separate community.

The algorithm iteratively improves modularity by moving nodes between communities. In each iteration, it evaluates the modularity gained by moving a node to a neighboring community or keeping it in its current community. The node is moved to the community which results in the maximum modularity gain. If no gain is achieved, the node remains in its current community. This process is repeated until no further improvement in modularity can be achieved.

3. Aggregation Phase:

Once the first phase of the Louvain Algorithm is completed, the network is transformed into a new network where communities become nodes. The weights of the edges between the new nodes represent the sum of the weights of the original network's edges between the corresponding communities. The algorithm then proceeds to the next phase, treating the newly formed communities as nodes and optimizing modularity in this aggregated network. The aggregation and optimization phases are repeated until a desired modularity level is reached or no further improvement can be made.

4. Output:

The Louvain Algorithm produces a partition of the network into communities, where each node belongs to one specific community. The resulting communities are non-overlapping, and the algorithm aims to maximize the modularity of the network. The Louvain Modularity Algorithm is known for its computational efficiency, making it suitable for large-scale networks. It has been successfully applied in various fields, including social network analysis, biological networks, citation networks, and recommendation systems. Its ability to detect communities with high modularity has made it a widely adopted algorithm in community detection tasks. It's worth noting that there can be variations and improvements to the Louvain Algorithm, such as resolution modifications and multi-level extensions, which further enhance its performance and accuracy in identifying communities within networks.

**3. Girvan-Newman Algorithm:**

The Girvan-Newman Algorithm is a hierarchical community detection algorithm that iteratively removes edges from a network based on their betweenness centrality to uncover the underlying community structure. It was proposed by Michelle Girvan and Mark Newman in 2002.

Here is an overview of the Girvan-Newman Algorithm:

1. Betweenness Centrality:

The algorithm utilizes the concept of betweenness centrality, which quantifies the importance of an edge in facilitating communication between nodes. Betweenness

centrality measures the number of shortest paths in the network that pass through a given edge. Edges with high betweenness centrality are considered as crucial for connecting different parts of the network.

2. Algorithm Steps:

Initially, the algorithm calculates the betweenness centrality for all edges in the network. It then removes the edge(s) with the highest betweenness centrality, thereby disconnecting the network into two or more separate components. The betweenness centrality calculation is recalculated for the updated network, and the process of edge removal is repeated. The algorithm continues removing edges iteratively until the desired number of communities or a stopping criterion is reached.

3. Community Hierarchy:

The Girvan-Newman Algorithm produces a hierarchical structure of communities. As edges are removed, the network breaks into smaller components, representing communities at different levels of granularity. The resulting dendrogram represents the nested community structure, where the top-level communities represent the largest clusters, and the bottom-level communities represent the smallest clusters.

4. Modularity Calculation:

Modularity is often used to evaluate the quality of the community structure produced by the Girvan-Newman Algorithm. After each edge removal, the modularity of the remaining network is calculated to measure the quality of the communities discovered so far. Higher modularity values indicate a stronger division of the network into communities.
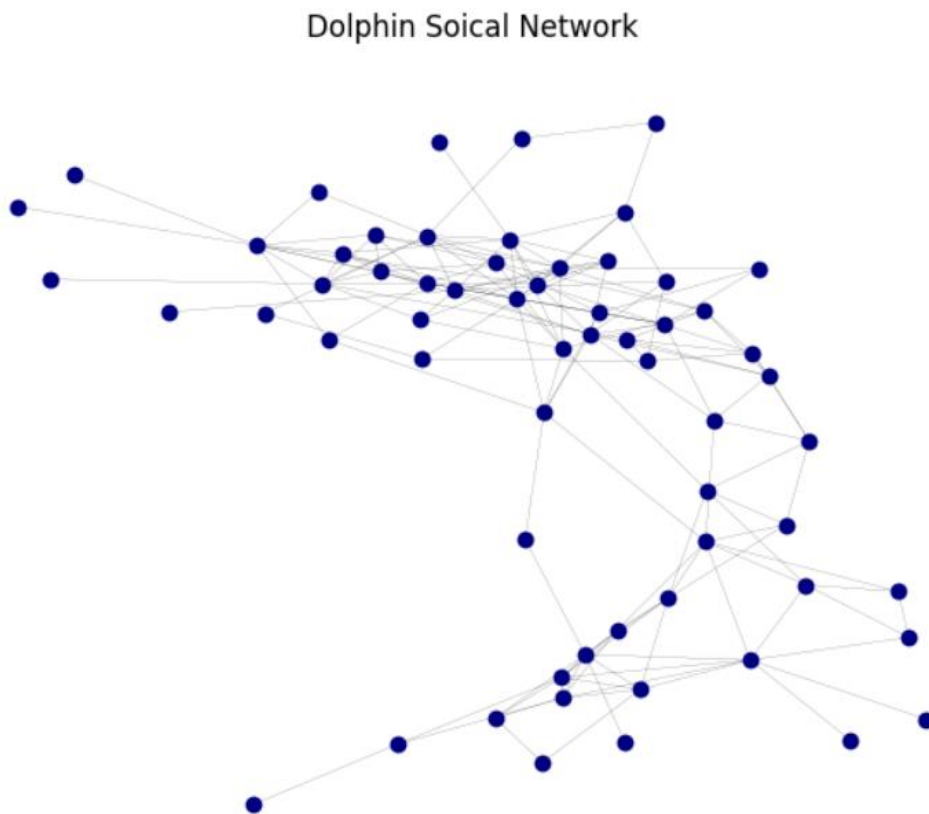

The Girvan-Newman Algorithm provides a systematic approach to community detection by progressively breaking down a network into smaller components based on the importance of edges in maintaining network connectivity. It is capable of detecting both hierarchical and overlapping communities. However, it is important to note that the algorithm can be computationally expensive, especially for large networks, as it requires calculating the betweenness centrality for all edges. Various optimization techniques, such as using approximate methods or heuristics, have been proposed to overcome this limitation.

The Girvan-Newman Algorithm has been widely applied in various domains, including social networks, biological networks, citation networks, and web networks, to uncover the community structure and gain insights into the organization and functioning of complex systems.
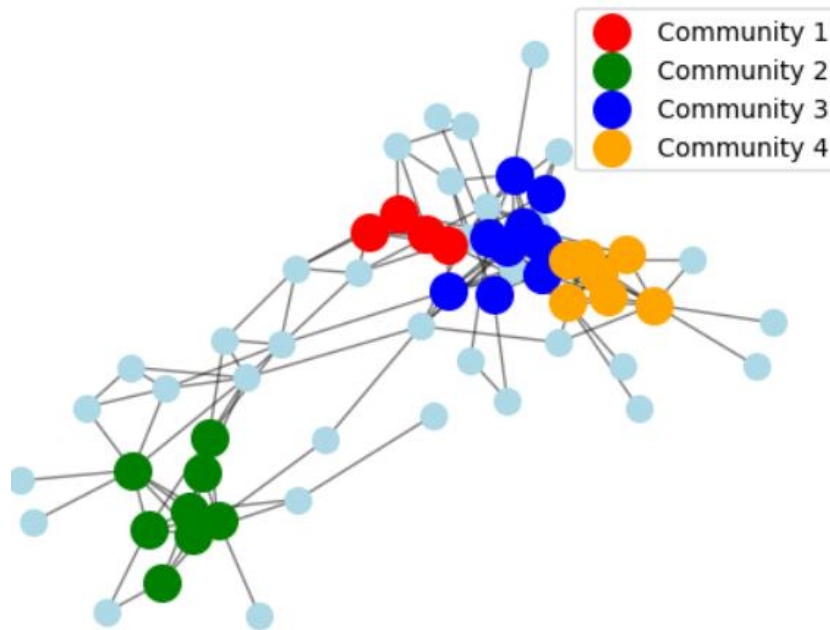
# 5. KEY RESULTS & CONCLUSION:
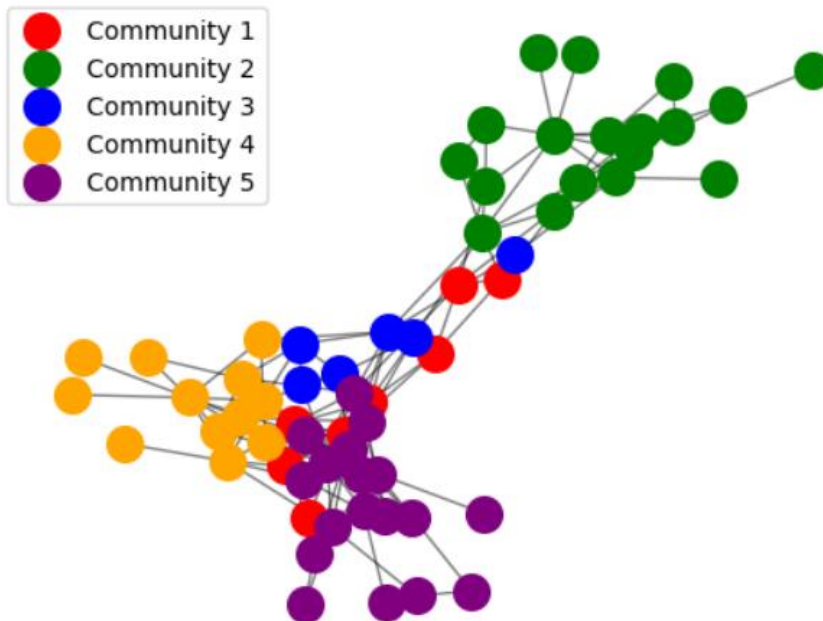
**1. Dolphin Social Network:**

Dolphin Soical Network



**2. Clique Percolation Method (CPM):**
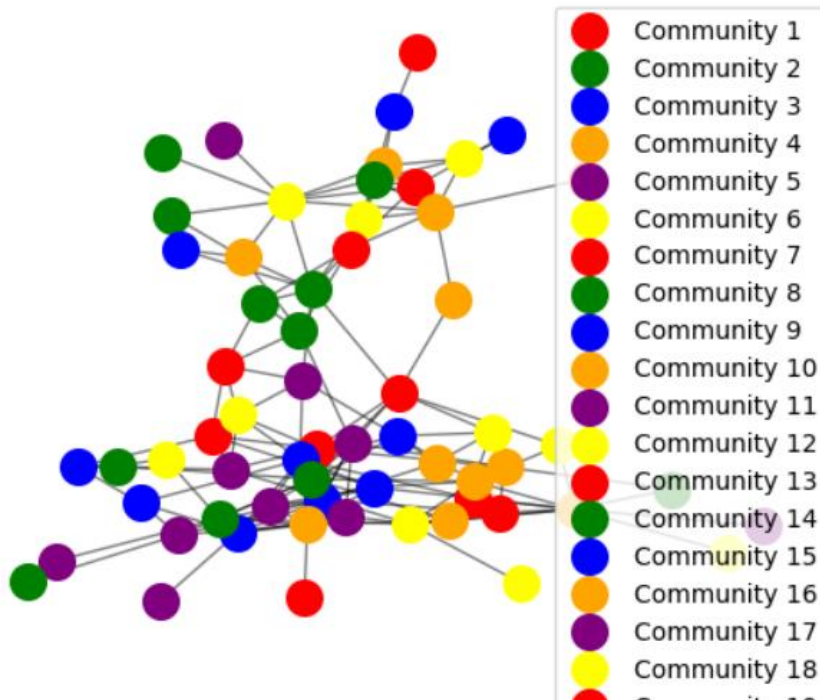
## Community Detection using CPM Algorithm



# 3. Louvain Modularity Algorithm:

## Community Detection using Louvain Algorithm



# 4. Girvan-Newman Algorithm:

Community Detection using Girvan-Newman Algorithm



| | |
|---|---|
| 🔴 | Community 1 |
| 🟢 | Community 2 |
| 🔵 | Community 3 |
| 🟠 | Community 4 |
| 🟣 | Community 5 |
| 🟡 | Community 6 |
| 🔴 | Community 7 |
| 🟢 | Community 8 |
| 🔵 | Community 9 |
| 🟠 | Community 10 |
| 🟣 | Community 11 |
| 🟡 | Community 12 |
| 🔴 | Community 13 |
| 🟢 | Community 14 |
| 🔵 | Community 15 |
| 🟠 | Community 16 |
| 🟣 | Community 17 |
| 🟡 | Community 18 |

After applying community detection algorithms such as Clique Percolation Method (CPM), Louvain Modularity Algorithm, and Girvan-Newman Algorithm to a social network, we can draw the following conclusions:

**1. Clique Percolation Method (CPM):**

- CPM identifies communities based on the concept of k-cliques, where nodes within a k-clique are considered to belong to the same community.
- CPM allows for overlapping communities, where nodes can be part of multiple communities simultaneously.
- The algorithm provides a systematic approach to identifying communities by considering the presence of shared k-cliques.

**2. Louvain Modularity Algorithm:**

- The Louvain Modularity Algorithm is known for its efficiency in identifying communities with high modularity in large-scale networks.
- It optimizes the modularity metric, which measures the density of edges within communities compared to a null model.

- The algorithm iteratively moves nodes between communities to improve the modularity, resulting in a partition of the network into communities.
- Louvain Modularity Algorithm produces non-overlapping communities and is widely used in various domains due to its computational efficiency.

**3. Girvan-Newman Algorithm:**

- The Girvan-Newman Algorithm is a hierarchical community detection algorithm that iteratively removes edges based on their betweenness centrality.
- By progressively removing edges with high betweenness centrality, the algorithm uncovers the community structure of the network.
- Girvan-Newman Algorithm produces a hierarchical community hierarchy, where larger communities are represented at higher levels and smaller communities at lower levels.
- The algorithm provides insights into the connectivity patterns within the network and helps identify communities at different levels of granularity.
- Overall, these community detection algorithms offer different approaches to uncovering the community structure in social networks. CPM focuses on the presence of shared k-cliques, Louvain Modularity Algorithm optimizes modularity to identify non-overlapping communities efficiently, and Girvan-Newman Algorithm utilizes betweenness centrality to reveal hierarchical community structures. The choice of algorithm depends on the specific characteristics of the network and the goals of the analysis. By applying these algorithms, we can gain a better understanding of the organization, dynamics, and functional relationships within the social network.

# 6. FUTURE SCOPE:

The field of community detection algorithms in social network analysis continues to evolve, and there are several potential future directions and areas of improvement:

1. Scalability: One major challenge is handling large-scale networks with millions or billions of nodes and edges. Future research can focus on developing algorithms that are more scalable and efficient, enabling community detection on massive networks in a reasonable amount of time.

2. Overlapping Communities: Many real-world networks exhibit overlapping community structures, where nodes can belong to multiple communities simultaneously. Enhancing algorithms to effectively identify and represent overlapping communities is an active area of research. Future work can focus on developing algorithms that can handle overlapping communities more accurately and efficiently.

3. Dynamic Networks: Real-world networks often evolve over time, where nodes and edges can be added or removed. Future community detection algorithms should consider the dynamic nature of networks and adapt to changes in community structures over time. Dynamic community detection can provide insights into the evolution and behavior of communities in dynamic social networks.

4. Multilayer Networks: Networks with multiple types of relationships or layers, such as social networks with different types of interactions (e.g., friendship, collaboration), pose unique challenges for community detection. Future research can focus on developing algorithms that can effectively handle multilayer networks and capture the interactions between different layers to identify communities more accurately.

5. Evaluation Metrics: Developing robust evaluation metrics for community detection algorithms remains an important area of research. Current metrics, such as modularity, have limitations and may not fully capture the quality and meaningfulness of communities. Future work can focus on developing new evaluation metrics that consider various aspects, such as the internal cohesion of communities, external separation, and functional relevance.

6. Domain-Specific Applications: Community detection algorithms can be applied to various domains, such as social media analysis, biological networks, recommendation systems, and cybersecurity. Future research can explore domain-specific applications and develop specialized algorithms that leverage domain-specific knowledge and characteristics to improve the accuracy and effectiveness of community detection.

7. Combination of Algorithms: Ensemble methods or hybrid approaches that combine multiple community detection algorithms can potentially improve the robustness and accuracy of community detection. Future work can focus on developing techniques to integrate different algorithms and leverage their strengths to achieve better community detection results.

In conclusion, the future of community detection algorithms in social network analysis lies in addressing scalability challenges, handling overlapping and dynamic networks, incorporating multilayer network analysis, improving evaluation metrics, exploring domain-specific applications, and exploring ensemble and hybrid approaches. These advancements will contribute to more accurate and meaningful community detection and provide deeper insights into the structure and dynamics of social networks.