

COMMUNITY DETECTION IN SOCIAL NETWORK USING GRAPH CLUSTERING METHODS

¹SURYANENI PRIYANKA, ²SAI RAMA KRISHNA

¹M.Tech student, Dept. of CSE, Kakatiya Institute of Technology& Science, Warangal, TS, srk.cse@kitsw.ac.in

²Assistant Professor, Dept. of CSE, Kakatiya Institute of Technology& Science, Warangal, TS, M19SE020@kitsw.ac.in

Abstract: Community detection is a major analysis area in social media analysis where we identify the social network construction. Detecting communities is of main interest in sociology, computer science, biology, and methods, where systems are usually described as graphs. Community detection aims at discovery clusters as subgraphs within a specified network. A community is then a cluster where many edges link nodes of the same group and few edges link nodes of different clusters. With the democratization of the internet, communicating and sharing information is more manageable than ever. Community detection is a solution to understanding the structure of complex networks and finally extracting useful information from them. In Facebook, the existence of communities (groups) is a critical question; thus, many researchers focus on potential communities by using techniques like data mining and web mining. In this paper, the community detection for a Facebook social network is presenting. It has developed in network science to find groups within complex systems depicted on a graph. The proposed model divided into various phases such as a) sub-graph discovery, b) vertex clustering, c) community quality optimization, d) divisive, and e) model-based. For community detection defining the consistency between social particles, Social media applied Social Balance, Social status theory, Social correlation theory, and finally applying K-means clustering over facebook data set. The Experimental conducted on Face book social media dataset with multiple edges and nodes. The experimental results shows that the proposed model gets higher accuracy in community detection compared with state of the art methods.

Keywords: community detection, social media, large scale networks, K-means, Facebook,

I. INTRODUCTION

Facebook and Twitter are social characters similar to people or companies that relate to each other through social interactions. It describes a dynamic social structure through a set of nodes and hyperlinks. The social networking assessment, which relies mainly on graph theories and social analysis, aims to look at the unique aspects of these networks. The most important elements are revealing the community, the identity of the influencing actors, and the shape and prediction of the development of the networks.

Social networks (among other complex networks) usually show the shape of a community (which is sometimes called an assembly). A network is said to exhibit such a structure if the network vertices can be divided into units of overlapping or separate heads. The number of inner edges exceeds the number of outer edges by a reasonable amount. The inner side is the edge that connects two heads belonging to the same network, while the outer part is the part that connects the heads of different communities. A social media is created for the individual through his private interactions and relationships with other individuals within the community. Social networks represent the social connections between individuals and their mobility. With the rapid growth of the network, there is a huge increase in user interaction online. For example, many social networking sites have appeared, for example, Facebook, Twitter, and others, to facilitate consumer interaction. As the number of interactions doubled, it became difficult to track those communications. Humans tend to associate with people of similar tastes and hobbies. The user-friendly social networks allow people to improve their social lifestyle in an unprecedented way. It isn't easy to satisfy friends in the world of the body. Still, it is much easier to find friends who will go along with similar hobbies. These real global social networks contain exciting patterns and homes that can analyze for various useful functions. A standard approach to community detection is to think of society as a static view in which all nodes and connections within a network remain unchanged at some point in

the test. Modern awareness research also looks at the evolution of society as most social networks tend to conform over the years by adding and removing nodes and links. As a result, groups within the network can grow or contract, and the participants in them can move from one organization to another over the years. In community detection to represent the graph, Social networking Site (SNS) like Facebook, Twitter can be modelled as a graph $G= (V, E)$, where V is a set of nodes and E is a set of edges that denote the communication among the nodes as shown in Fig.1.

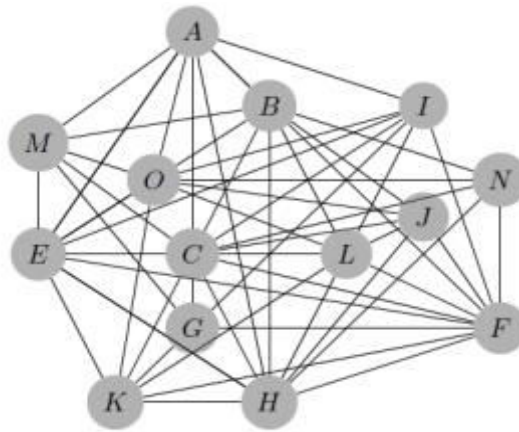


Fig.1 Social media network

Social media has the membership feature to show some form of community. Suppose the community vertices can be divided into units of both discrete and overlapping vertices. The edges range within the edge diversity group exceeds any two units using a reasonable amount. In that case, we say that the community presents the network structure. Networks that regularly display a network topology may display a hierarchical network structure as well. Most research on network evolution uses topological characteristics to select the updated elements of the community and symbolize the type of modification, such as network reduction, development, segmentation, and merging. However, the current panel has focused on network evolution/discovery by relying entirely on individuals' behavior in terms of activities originating within the community rather than just looking at links and network density.

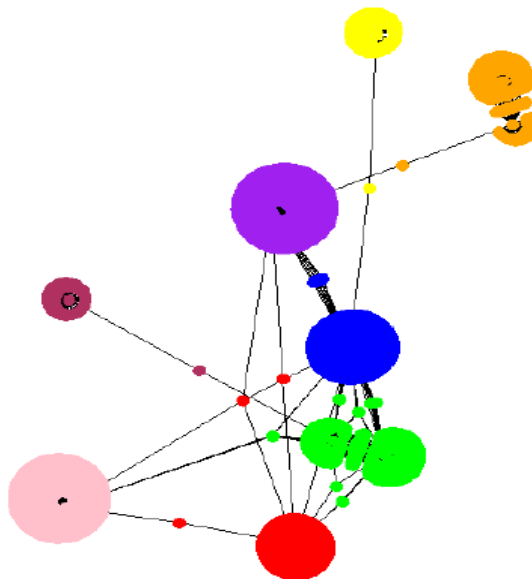


Fig.2 Facebook NIPS Social Network

Figure 2 shows the example network graph for PGP Network Dataset. This dataset is the graph of a component of a network, included the users of the Pretty-Good-Privacy algorithm for privacy data exchange. PGP network contains 10680 vertices and 24316 edges.

II. RELATED WORK

Early community detection algorithms were mostly based on graph segmentation theory and hierarchical clustering in sociology.

Guangxia et al. [2020] Social networks are a type of community made up of many nodes according to unique types of relationships. Community discovery can help people understand the topology of a community and identify important groupings. A unique approach to community detection algorithms has been proposed. However, many of these algorithms are suitable for non-weighted networks because they ignore the social housing of the hyperlinks between nodes. Meanwhile, the consequences of spotting them can also give rise to large clusters. Therefore, this post proposes a weighted and close societal stabilization method for social networks that relies entirely on LFM. This WSLFM algorithm uses a weighted technique that relies on the social attributes of the association and the degree of common neighboring nodes to update the health characteristic. It also introduces the concept of stability to govern the expansion of local communities. The proposed WSLFM algorithm has been examined and demonstrated its robustness in an artificial community and a truly international network to locate smaller, more important clusters.

Vasima et al. [2020] with the significant increase in consumer data, their use of research works has become very popular. This provides an upward push to discover society as a growing discipline in social media exploration. Also, the researchers became interested and began working in this area. But there are some limitations in this area, such as scalability and the first class of the network. Some search algorithms have completed the work in this area, and they are also efficient in terms of scalability and accuracy. We contrasted several algorithms with the one proposed in Twitter's Social Facts. With the help of our results, we showed that our set of rules is more efficient at evaluating legacy algorithms.

Muhammad Aqib et al. [2018] the advanced technological knowledge of networks has seen a major advance in the modeling of complex real international systems. One of the most important capabilities of these networks is the existence of a community structure. In recent years, several community detection algorithms have been proposed to detect the structural houses and dynamic behaviors of networks. In this insight, we strive for a contemporary survey of network detection techniques and packages in the names of various domains of real existence. In addition to highlighting the strengths and weaknesses of each network discovery approach, various elements are mentioned for comparing the performance of the algorithm and testing them against benchmarks. Challenges for network discovery algorithms, open issues, and future trends related to community discovery are also discussed. The main objective of this article is to provide a review of winning community discovery algorithms that range from traditional algorithms to next-generation nested network detection algorithms. Algorithms based on dimensional minimization strategies also focus on, including non-negative matrix factorization (NMF) and principal component analysis (PCA). This overview will serve as an up-to-date archive on the evolution of the community's discovery and bundling of capabilities on the many real-world network domain names.

Xiaolonget et al. [2017] Community detection algorithms are essential for identifying male or female information from complex networks. Compared to traditional community detection algorithms, which are generally recognized in unguided networks, our algorithm focuses on directed networks, including WeChat's moment of courtship and follower relationship network in the Sina Micro-Blog. To address the downsides, including the low runtime performance and high-precision bias that target community detection algorithms consistently enjoy today, we recommend a completely new method that relies primarily on the triple topology of the network base and is designed based on neighborhood information transmission technology. It completely collides with groups in the target networks. Based on the concept of vector in probability diagrams and the dynamic information transfer gain (ITG) of heads in directed networks, we recommend the new Integrated Technologies Group (ITG) approach and the corresponding optimized objective function for comparing the sub-section in a set of

community discovery rules. Then we combined the Integrated Technologies Group (ITG) and the target function to create the new network discovery algorithm, weighted network aggregation led by Integrated Technologies Group (ITG) for directed networks. With full-size experiments using artificial network information units and large sets of real global community statistics drawn from online social networks, our algorithm has proven to be healthier and faster on target networks than many popular community discovery technologies. And traditional ones, including FastGN, Neighborhood Improvement Technology Records, and the info map.

Girvan et al. [2016] we have revealed a specific equivalence between two network discovery strategies widely used in networks, the modular maximization approach in its generalized form, which includes a decision parameter that controls the dimensions of the identified clusters, and the maximum likelihood method applied to the unique case of randomness. The cluster model is called the implanted partition model, where it is assumed that all groups in the network have statistically similar properties. Among other things, this parity provides a mathematical derivation of a modular property, clarifies the conditions and assumptions for its use, and provides explicit components of the price most appropriate to the coefficient of accuracy.

Zhang et al [2015] many networks can advantageously decompose into a dense medium as well as a loosely bound peripheral periphery. We recommend an algorithm to perform this type of analysis on experimental network records using statistical inference techniques. Our technology adapts a generative model of the shape of the outer edge of the core to the detected information using a combination of a predictive maximization algorithm to calculate release parameters and an idea propagation algorithm to calculate the decomposition itself. We localized the technology to be effective, scaled without difficulty into networks of a million or more nodes, and tested it in very few networks, including real international examples and benchmarks built by laptops, which define correctly recognized cores. - Circumferential structure with low fault rate. Moreover, we show that this technique is immune to the discoverability transmission within the relevant community discovery problem, which avoids discovery of the community structure while this structure is very weak. There is no such transition to the shape of the outer edge of the core, which could be detected, albeit with some statistical errors, regardless of its susceptibility.

Dominguez et al. [2014] Discovering communities has emerged as one of the most relevant topics in graph mining, mainly due to its programs in social networking or membership assessment. Various algorithms for revealing the community have been proposed over the past decade, nearly annoying the unusual viewpoints. However, current algorithms rely mainly on complex and fascinating computations, which makes them incorrect for large graphics with tens of millions of vertices and edges and those typically defined in the real world. This document recommends a new, separate network detection algorithm known as Scalable Community Discovery (SCD). By combining unique techniques, SCD wall graph by weighted community clustering (WCC), a newly proposed community finding scale based on triangulation evaluation. Using simple graphics with overlapping communities of terrain reality, we show that SCD outperforms the modern state of technical proposals (even those aimed at creating overlapping combinations) in terms of first-class and performance. SCD offers the fastest algorithmic speed and pleasure in NMI syntax, and F1Score is the right extreme proposition for the art field. By exploiting parallelism with modern multimedia processors, we demonstrate that SCD can perform up to two important things faster than real-world responses, allowing us to process significantly long graphics at fast run times.

III. PROPOSED METHODOLOGY

Community detection in these large-scale social networks plays an essential role in analyzing network topologies and structures. Due to the large number of vertices and edges in a large-scale network and its complex shape, common traditional graph analysis procedures cannot perform the investigation, including analyzing layered structures and insight into cases. Affordable implementation. Therefore, over-productivity and correct community discovery algorithms are vital in locating community structures of capacity in a large, complex target community.

The proposed work divided into various phases of (a) sub-graph discovery, (b) vertex clustering, (c) community quality optimization, (d) divisive, and (e) model-based

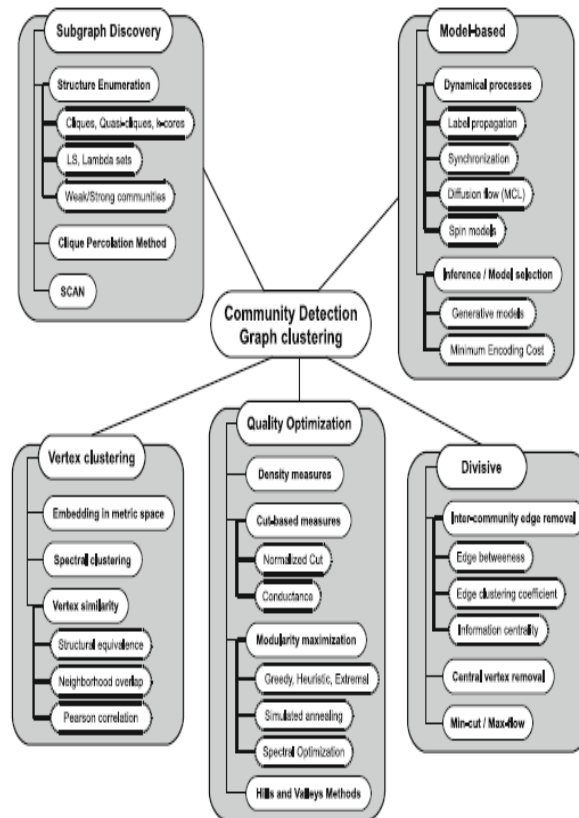


Fig.3 A classification of community detection and graph clustering methods

a) Sub graph discovery

Strategies for this elegance presuppose a determination of the structural properties that a sub-graph of a network must fulfill to be considered a community. Once this type of sub-graph diagram is private, the strategies contain an enumeration of these structures within the network below.

b) Vertex clustering

These techniques arise from traditional data clustering studies. A typical way to convert a graph vertex grouping problem into a problem that can be solved using traditional information aggregation methods (including k method and hierarchical clustering method) is with the help of including graph vertices in a vector region, where the actual distances between the vertices can be calculated. Another popular technique is to use graph spectrum to map graph headers with agents in a low-dimensional space, where the group shape is deeper.

c) Community quality optimization

Various strategies can be based on improving a certain degree of excellent society that relies primarily on graphics. Density and subshell plot measurements, including standard cut and conduction, were among the first measurements used to determine the excellent level of a few sections of the community in groups.

d) Divisive

These strategies are based on the identity of community factors (edges and vertices) that are placed between communities. For example, the basic algorithm using Girvan and Newman (2002) gradually removes population boundaries based on some measures of lateral mediation until the groups become separate components of the graph. Several measures of correlation between faces have been developed, for example, the correlation between faces, random walk, and recent wave, in addition to data centrality and threshold aggregation coefficient.

e) Model-based

It is a huge class of unique styles, the dynamic one that takes the neighborhood out of the community, shows its communities well, or does not forget a basic version of the statistical nature that a network division can create.

f) Performance comparison

When evaluating the overall performance of community discovery methods, there are two primary components one would like to consider: (a) computational complexity, and (b) technology needs in terms of memory are more important.

COMMUNITY DETECTION ON SOCIAL NETWORKING

Social networking SN can be denoted as a graph $G(P, R, W)$. Where P is set of peoples (vertices) belong to Social network, R is a set of relationship among two elements of P and $W: p \times p \rightarrow R$ is a function which allocates a weight to a couple (P_i, P_j) of vertices P_i and P_j , for instance if $W: P_i \times P_j \rightarrow 1$ then their having an connection between P_i and P_j . Whereas if $W: P_i \times P_j \rightarrow 0$ then there is no communication among P_i and P_j . Social networking websites do not distribute real dataset social network. Before publishing people's data, social media owners anonymized social media statistics using traditional anonymous method (like; k-anonymity, t-closeness, and i-diversity). Anonymized social networks data can be represented with the adjacency matrix $AP \times P$ and value of A_{ij} determine the type of network. If $A_{ij} = A_{jii}AP \times P$ is symmetric matrix then SN sites is undirected network

Community collect similar social atoms with no geographic boundaries and have a similar view on social, political, monetary, and global issues on a social media platform. The goal of network discovery is to discover the organization of headers (sub-graphs) that have a high density of hyperlinks within the organization and less density of hyperlinks outside the organization's gates.

Based on the above discussed phases the connected graphs for community detection determining the consistency among social atoms, Social media applied Social Balance, Social status theory, Social correlation theory, and finally applying K-means clustering over social media data set.

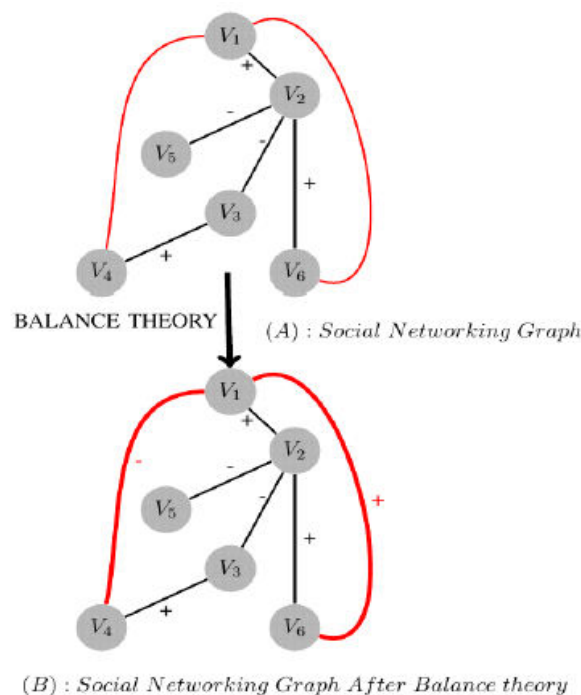


Fig.4 Social balance theory

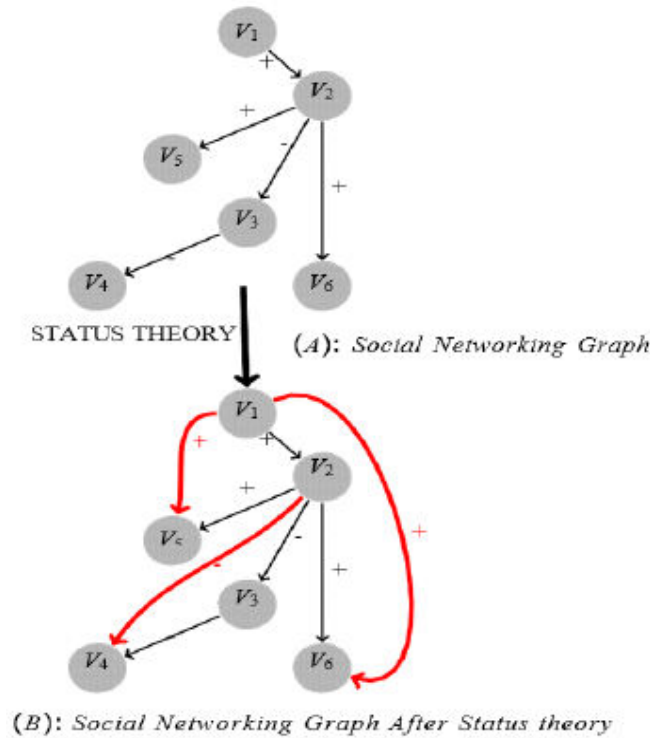


Fig.5 Social status theory

For example, don't forget the social graph shown in Figure 4 (A), where the red colored Nodes of the Republican political celebration and the inexperienced shadow nodes are politically neutral. Due to the principle of title correlation, sending the scarlet knot and its permanent messages get a color knot without influence experience to become a follower of the republican political celebration as shown in Figure 4 (B). Whereas homologous, the organization of the name of social atoms (nodes) in their coloring notation, as shown in Fig. 4 (A), while confusing world environments make people alike. Two people who live in the same city are more likely to become friends than two random people.

Community detection algorithm

Input: $SM G = (V, E)$

% De – Anonymization of social network graph

{

% Apply social status theory

{

Step1: *Ordering the vertex $(V_1, V_2, \dots, V_n$ in accending*

order in according to their degree

Step 2: *Ordering the vertex on behalf of their*

degree such that for every edge $X \rightarrow Y, X$

have higher status then Y

}

% Apply Social balance Theory

{

Step 3: *Placing an edge E_i between X and Z such that*

there exist an edge $X \rightarrow Y$ and $Y \rightarrow Z$

}

% Apply Social Correlation Theory

{

Step 4: *If there is a Edge Between $X \rightarrow Y$ (Status of X*

Is higher than Y) Then Y is Influence By X

}

}

Construction of adjacency matrix with vertex degree

vector

{

Step 5: *Construct an matrix $[P]_{N \times N}$ Where $N = \sum V$*

Step 6: *Insert value in $[P]_{N \times N}$ such that*

$V_1 \rightarrow V_2 \begin{cases} \text{Yes Place 1 at } [P]_{V_1 V_2} \\ \text{No: Place 0 at } [P]_{V_1 V_2} \end{cases}$

} % Vertex similarity

Step 7: *Construct one dimensional Matrix $[V]_{1 \times N}$*

Step 8: *Insert value in $[V]_{1 \times N}$ such that*

$$[V]_{1 \times N} = \sum_{j=1}^{j=N} P_{1 \times N}$$

}

*% Apply K – means clustering over $MaxV_{1*j}$ vertex*

{

Step 9: *K – means ($[P]_{N \times N}$ by choosing $MaxV_{1*j}$*

as initial point

}

IV. RESULTS AND DISCUSSIONS

The proposed model experimental taken on the Facebook dataset with multiple edges and nodes.

Experimental setup

Implementation of our algorithm is done with Python 3.5. For visualization, Graph- Stream (<http://graphstream-project.org/>), which is a java library to model and analyze dynamic graphs is used. For the experiments for K-mean clustering algorithm. Jypitor is a tool kit for community detection in networks is used. Hardware specifications of the computer used for experiments are following: Windows 10 Operating System with 64 bit, 6 GB RAM as memory of which 4000m is allocated to Python Heap Space, Intel(R) Core(TM) i3-2520M 2.50 GHz CPU as processor.

Community detection for 6 nodes and 8 edges with average degree 2.6667 as shown below

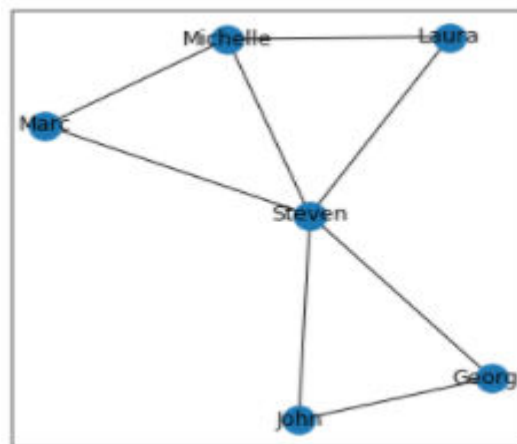


Fig.6 Graph with 6 nodes and 8 edges.

Figure.6 represent the connected graph between 6 nodes and 8 edges

The below figure 7 shows the connected graph with number of nodes 4039 and edges 88234 with average degree is 43.6910.

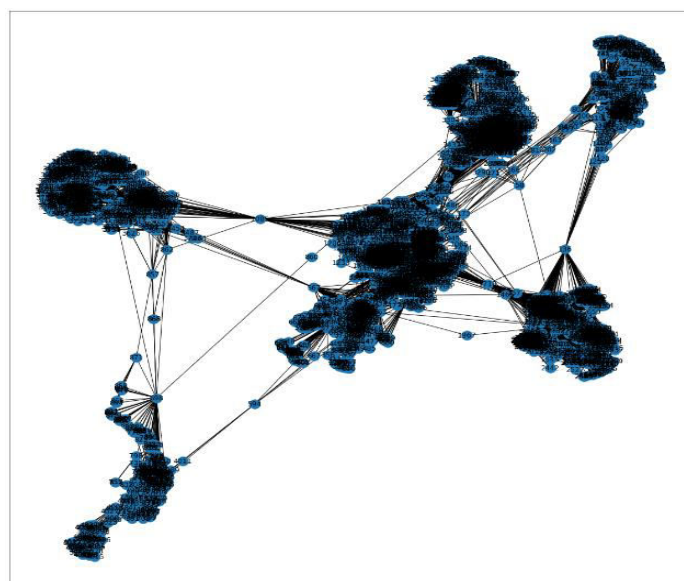


Fig.7 Graph with nodes 4039 and edges 88234

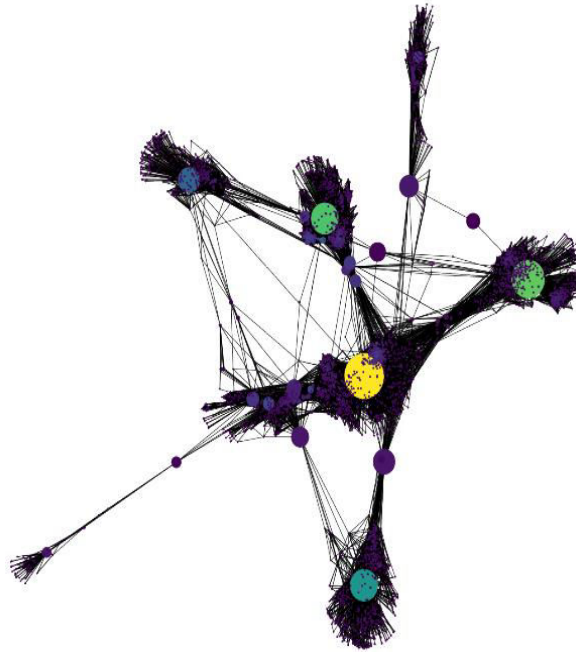


Fig.8 Graph with multiple nodes and multiple edges

Figure.8 shows the connected graph with multiple nodes and multiple edges on facebook related dataset.

V. CONCLUSION

Community detection has been tirelessly studied for these years, and an especially wide range of effective processes has been proposed. This observation recognizes the basic problem of detecting non-overlapping networks, and the objectives are to find the specific group (community) to which each node in the graph belongs. In this paper, the community detection for Facebook social networks is presented. It has developed in network science to find groups within complex systems depicted on a graph. Social media applied Social Balance, Social status theory, Social correlation theory, and finally applying K-means clustering. The proposed model experimental conducted on the 'facebook_combined' dataset with multiple edges and nodes and the connected graphs for various edges and nodes presented.

REFERENCES

1. P. J. Mucha and A. L. Traud, 2011, "Social structure of facebook networks,"
2. Z. Kobti and Zadeh P. M, 2015, "A multi-population cultural algorithm for community detection in social networks", pp. 342 – 349.
3. X. Niu and Wu, C. Q, 2017, "A label-based evolutionary computing approach to dynamic community detection," pp. 110 – 122.
4. Guangxia Xu and Liu Yanbing, 2020, "A Community Detection Method Based on Local Optimization in Social Networks", IEEE, pp.42-48.
5. Vasima Khan and Manoriya Manish, 2020, "Influence Based Community Detection Over Social Media", pp.
6. Muhammad Aqib Javed and Junaid Qadir, 2018, "Community detection in networks: A multidisciplinary review".
7. D. Xiaolong, Zhai Jiayu, 2017, "Efficient Vector Influence Clustering Coefficient Based Directed Community Detection Method", pp.17106-17116.
8. M. Girvan and M. E. J. Newman, 2016, "Community detection in networks: Modularity optimization and maximum likelihood are equivalent," pp. 052315.

9. X. Zhang and T. Martin,2015, ``Identification of core-periphery structure in networks," p. 032803.
10. D. Dominguez-Sal and A. Prat-Pérez, 2014, ``High quality, scalable and parallel community detection for large real graphs," pp. 225-236.