

# Facebook Dataset Analysis

Varang Pratap Singh

Data Science And Artificial Intelligence

Indian Institute of Information Technology, Dharwad  
Graphs And Social Network Project

---

**Abstract** — *The project's objective is to analyze the Facebook dataset to identify user interactions, influential nodes, and community structures. Utilizing network analysis methodologies, we constructed and analyzed the network employing Python libraries such as NetworkX. Visualization techniques, including network metrics, were employed to present our findings. Centrality metrics and community detection algorithms unveiled user influence and information dissemination within social clusters. The analysis demonstrated small-world phenomena of Facebook and existence of clustered communities facilitating information flow. This underscores the efficacy of network analysis in uncovering latent structures and patterns within digital social networks, providing valuable insights for social research and online community management.*

---

## I. INTRODUCTION

Social network data on platforms like Facebook provides researchers and data scientists with unique opportunities to study human behavior, interactions, and social connectivity. This data captures relationships, communities, and influence levels, making social networks crucial for understanding individual interactions within broader societal contexts and the spread of information and trends.

The report analyzes Facebook's social network using a dataset of 4,039 nodes and 88,234 edges, representing friendships between users. Each user is a node, and friendships are edges, enabling social network analysis techniques to examine connectivity, community structures, and influence dynamics.

The dataset shows strong connectivity, supporting the small-world phenomenon. High clustering coefficient indicates tight-knit communities. We analyzed the network's centrality, community structure, and clustering to reveal users' roles and information spread.

### **Centrality Measures**

One of the primary objectives is to identify influential users—those playing central roles in connecting and influencing others within the network. Centrality measures are metrics assessing the position and influence of nodes based on their connectivity patterns. Several centrality measures were used to identify key players:

#### **- Degree Centrality**

Counts the number of direct connections a node has. Nodes with a high degree centrality are considered “hubs” or highly connected individuals who can significantly impact their community. In Facebook dataset, Node 107 emerged as the most connected users with a degree centrality of 0.2588, indicating direct relationships with large number of other users.

#### **- Betweenness Centrality**

Identifies nodes acting as intermediaries between different parts of a network. These are critical for information flow, serving as bridges connecting distant clusters. Node 107, which also held the highest betweenness centrality (0.4805), plays a key role in facilitating communication across the network, making it a potential influencer in bridging various social groups.

#### **- Closeness Centrality**

Evaluates how quickly a node can reach all other nodes in the network. Nodes with high closeness centrality disseminate information more efficiently, as they are closer to all other nodes in terms of path length. Node 107 also had the highest closeness centrality score (0.4597), marking it as a user who could efficiently spread information through the network.

#### **- Eigenvector Centrality**

Assesses a node's influence by considering not only its direct connections but also the connections of its neighbors. Nodes with high eigenvector centrality are connected to other well-connected nodes, thus amplifying their influence. In this dataset, Node 1912 exhibited the highest eigenvector centrality score (0.0954), indicating its importance through associations with other influential nodes.

#### **- PageRank**

Evaluates a user's importance based on the quality and quantity of their connections. Node 3437 had the highest PageRank score (0.0076), highlighting it as a critical figure within the network, due to its connections with other prominent nodes.

### ***Community Detection***

We also explored the community structure of the network to understand how users group. Communities within a social network represent groups of users with denser connections among themselves than with users outside their community. Detecting these communities is essential for understanding the social circles and shared interests binding users together.

#### **- Louvain Method**

Identifies communities by optimizing modularity, a measure of the strength of divisions within the network. Efficient for large networks and successful in revealing distinct social circles within the dataset, where clusters of closely related users represent friend groups or social affiliations.

#### **- Label Propagation**

Detects communities by allowing nodes to adopt their neighbors' labels iteratively. This helps in revealing natural groupings. Label propagation effectively highlights the social circles in the network, suggesting presence of subgroups formed around common interests or mutual friendships.

#### **- Girvan-Newman Algorithm**

Finds communities by iteratively removing edges with high betweenness centrality. This approach revealed more granular divisions within network by identifying links playing a critical role in connecting separate communities. This further confirmed presence of densely connected groups within the network, vital for understanding Facebook's community dynamics.

### ***Clustering Coefficient and K-Core Decomposition***

Measures the extent to which nodes tend to cluster together, offering insights into the network's local cohesiveness. A high global clustering coefficient of 0.605 was observed, indicating that Facebook users are likely to form tightly connected groups. To further explore the network's structure, a k-core decomposition was performed, iteratively removing nodes with fewer than k connections. We identified tightly connected subgroups within the network.

### ***Network Structure and Connectivity***

The network's overall structure reflects small-world phenomena, where the average path length is short, and the maximum diameter is 8. This indicates that any user can reach any other within a maximum of eight steps. This facilitates rapid information flow, supporting the idea that social networks are highly interconnected. The combination of dense clustering and short average path length aligns with the small-world phenomena, suggesting that users can communicate and share information quickly despite relatively low density (0.0108).

## II.

## METHODOLOGY

The methodology involves a series of analytical techniques designed to uncover structural and relational patterns within the dataset, focusing on centrality, community detection, clustering, and other key network properties. Here are the techniques used:

### ***Data Preprocessing and Network Construction***

The first step was importing and preparing the data using NetworkX, a Python library tailored for network analysis. NetworkX facilitates creation, manipulation, and analysis of graphs, making it useful for understanding complex social networks.

- Each user is treated as a node, and edges represent friendships, making this an undirected graph with bidirectional connections.
- Calculating network parameters to understand dataset's properties:
  1. Density measures the ratio of actual edges to possible edges, showing how connected the network is.
  2. Diameter indicates the greatest distance between any two nodes, helping us understand the network's spread and navigability.
  3. Clustering Coefficient assesses the tendency for nodes to cluster together, indicating how closely knit the users are in their social circles.

Using these calculations, we establish a structural overview, giving insight into the interconnectedness of the network and laying foundation for further analysis.

### ***Centrality Measures***

Essential for identifying most influential or strategically positioned nodes within a network. Different measures provide unique perspectives on a node's influence and connectivity:

#### **- Degree Centrality**

Counting the number of direct connections each node has, this highlights users with high connectivity. These users (hubs) likely have significant reach and impact within their communities.

#### **- Betweenness Centrality**

Captures the extent to which a node serves as a bridge, connecting otherwise distant nodes. High betweenness centrality is critical for information flow, as they control the shortest paths between various clusters.

#### **- Closeness Centrality**

Identifies users who can reach other nodes with minimal steps. High closeness centrality suggests users can efficiently disseminate information across the network, given their proximity to other nodes.

#### **- Eigenvector Centrality**

Considers a node's connections to other well-connected nodes. Emphasizes the influence of nodes connected to other influential users, enhancing their social importance within the network.

#### **- PageRank**

Evaluates a node's prominence based on quantity and quality of its connections. It identifies users whose connections are influential, making them critical for information flow and reach.

These help us pinpoint influential nodes, essential for understanding information spread and user importance within the network.

### ***Community Detection***

Uncovers groups of nodes with denser connections within themselves than rest of the network. Revealing internal social structures and circles within the dataset:

#### **- Louvain Method**

Aims to identify clusters by maximizing modularity, a measure of the density of connections within clusters versus between clusters. Efficient for large datasets, it helps isolate tightly knit social circles, representing groups with shared interests or affiliations.

### **- Label Propagation**

Based on the principle that nodes adopt the labels of neighbors iteratively until stable community structures emerge. Reveals organic communities, reflecting natural social groupings without predefined labels.

### **- Girvan-Newman Algorithm**

Uses edge betweenness centrality to split the network. By progressively removing edges with high betweenness, this isolates clusters, identifying divisions based on critical connections. The Girvan-Newman method is particularly insightful for visualizing core and peripheral community structures within a network.

These algorithms provide a multidimensional view of Facebook's social circles, offering insights into user groupings and how different communities within the network interact.

## ***Graph Traversal Techniques***

Employed to understand shortest paths between nodes, enhancing our knowledge of network navigability. Two primary techniques were used:

### **- Breadth-First Search (BFS)**

Explores all neighbors of a node before moving deeper. In social networks, BFS helps find the shortest paths in terms of connections, highlighting how information might spread through immediate neighbors.

### **- Dijkstra's Algorithm**

Calculates shortest path from given node to all other nodes, accounting for varying weights in edge costs if applied. Though not strictly necessary in an unweighted social network, it is useful when studying scenarios where paths of varying importance exist, enhancing understanding of reachability.

These provide insights into network accessibility and user proximity, offering practical implications for understanding information flow within a network.

## ***Visualization***

Crucial for interpreting and presenting network structures. Various techniques were used to create a graphical representation of the network and its structural properties:

### **- Node-Link Diagrams**

Illustrates the entire network, representing users as nodes and friendships as edges, allowing clear view of how individual users connect to form the larger network.

### **- Clustering Coefficient Distribution**

Shows distribution of clustering coefficients, indicating the tendency of users to form clusters. High clustering in certain parts of the network reflects strong community structures, offering insights into local cohesiveness.

### **- K-Core Subgraph Visualization**

By visualizing different k-core decompositions (subgraphs formed by recursively removing nodes with degrees lower than k), we identify dense cores of highly interconnected users. This helps understand community structures and core influencers.

These visualizations make abstract properties more tangible, allowing for easier identification of influential nodes, clusters, and connectivity patterns.

### III.

## RESULTS AND DISCUSSION

The analysis of the dataset reveals several significant patterns deepening our understanding of its structure, interactions, and community dynamics. The findings offer valuable insights into how users connect, influencing one another, and forming cohesive groups within the network.

### *Network Density and Connectivity*

One of the first metrics analyzed was the network density, which measured at 0.0108. This indicates that, despite the large number of users (4,039) and friendships (88,234), the network remains relatively sparse. Such sparsity suggests that while there are connections between many users, a significant number of users don't have many direct friends, leading to a network where many possible connections remain unutilized.

Despite this sparsity, the maximum diameter of the network is 8, suggesting small-world phenomena. In this, most nodes can be reached from any other node by a small number of steps, reflecting the phenomena where individuals are connected through a series of acquaintances. This is crucial for information dissemination, as it implies that news, trends, or other content spreads rapidly across the network. Short average path length allows for quick communication between users, making the network efficient for sharing information.

### *Centrality Analysis*

These provide a clear picture of influential users within network.

#### **- Degree Centrality**

Node 107 emerged as the most connected user with a degree centrality of 0.2588. This value indicates that Node 107 has numerous connections, establishing it as a significant hub within the network. Users with high centrality are pivotal in influencing their communities due to direct access to a larger number of peers, facilitating effective information flow.

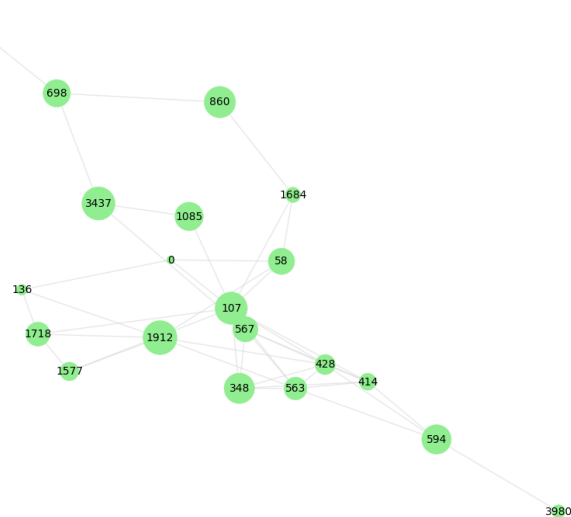
Top 20 Nodes by Degree Centrality



### - Betweenness Centrality

Node 107 recorded the highest betweenness centrality of 0.4805. This highlights its role as an intermediary between other users, indicating that Node 107 is crucial for connecting disparate parts of the network. Users with high betweenness centrality control information flow and serve as gatekeepers, especially in vital social networks.

Top 20 Nodes by Betweenness Centrality



### - Closeness Centrality

With a closeness centrality score of 0.4597, Node 107 is well-positioned to reach other users quickly. Node 107 can disseminate information efficiently across the network, reinforcing its influential status. Ability to quickly connect with other nodes enhances this user's role in facilitating communication within the social structure.

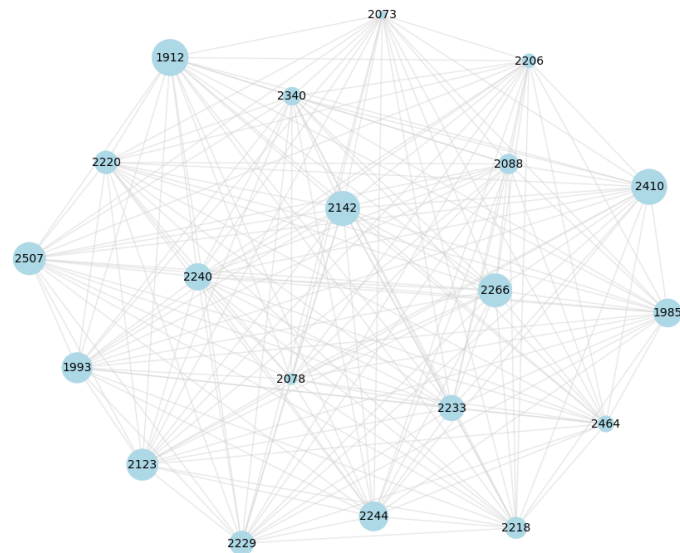
Top 20 Nodes by Closeness Centrality



### - Eigenvector Centrality

Node 1912 stood out with the highest eigenvector centrality of 0.0954. This indicates that Node 1912 is not just well-connected itself but is also connected to other influential nodes. This quality amplifies its importance in the network, as connections to influential users enhance its ability to spread information and trends.

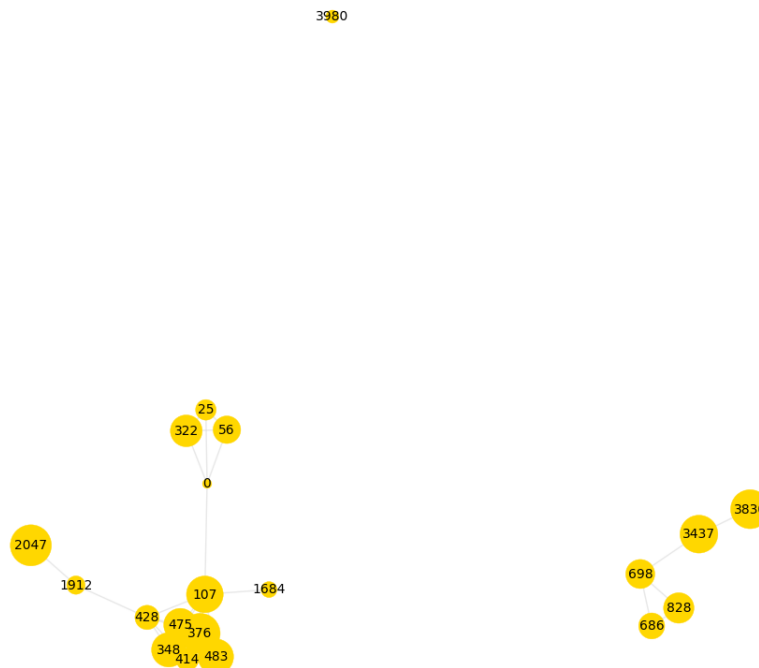
Top 20 Nodes by Eigenvector Centrality



### - PageRank

The score of 0.0076 for Node 3437 signifies its prominence based on both the quantity and quality of its connections. PageRank considers not just how many connections a user has but also how well-connected those connections are. Thus, a high PageRank indicates that Node 3437 is strategically positioned to affect social dynamics within the network, given its relationships with other influential users.

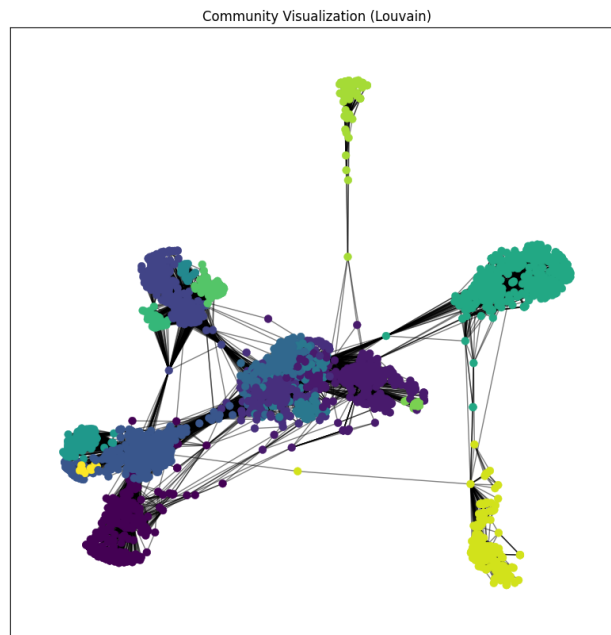
Top 20 Nodes by PageRank



## Community Structure

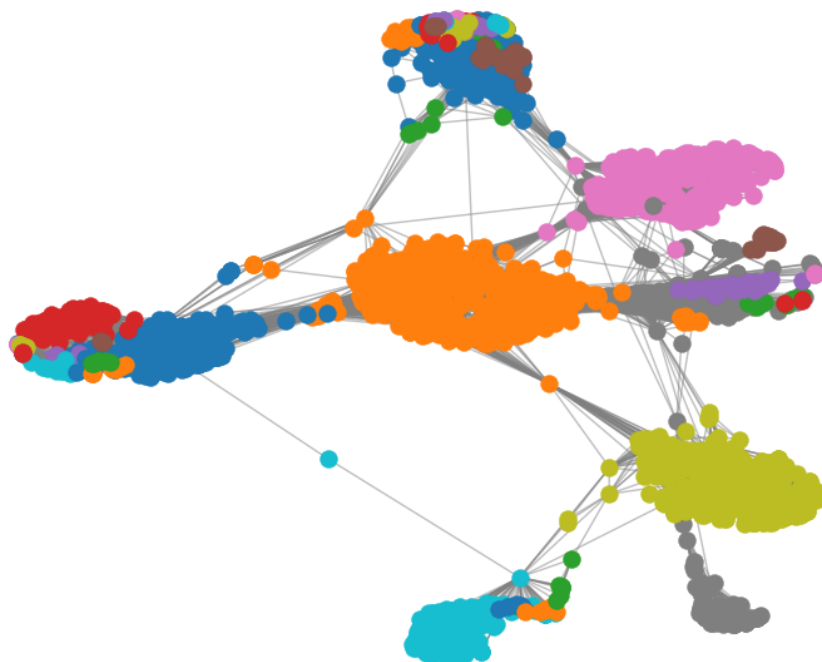
The analysis used community detection algorithms, including Louvain and Label Propagation, to identify tightly-knit communities reflecting social circles in the dataset. These clusters are crucial for understanding information spread and user interactions within groups.

The **Louvain** method maximizes modularity, revealing dense connections among users within communities while minimizing connections between them. This highlights the importance of community structures, as users within communities likely share interests or engage in similar activities, making targeted communication strategies more effective.



The **Label Propagation** algorithm showed how users naturally group into social circles by allowing nodes to adopt their neighbors' labels iteratively. It revealed insights into social dynamics.

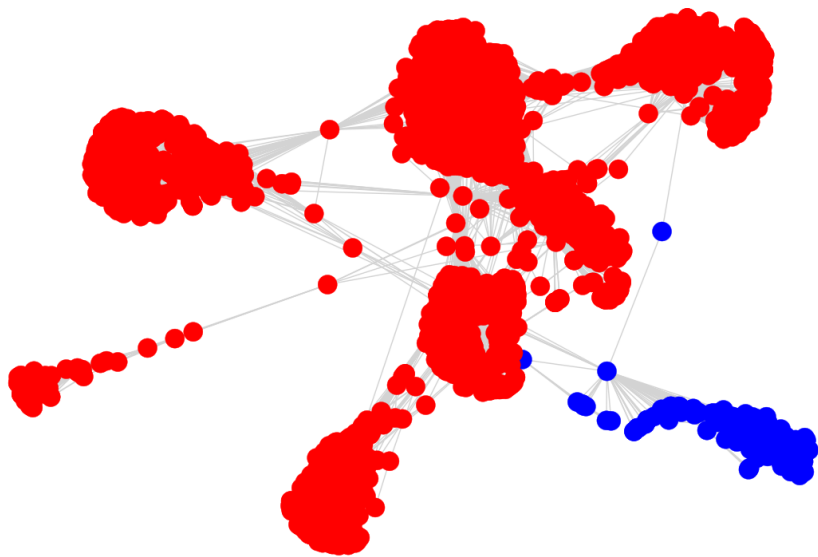
Community Visualization (Label Propagation)





The **Girvan-Newman** algorithm detected high-traffic links, highlighting critical connections between groups. It also identified bridges between communities, which facilitate cross-community information flow and interactions. Understanding these inter-community connections aids in effective communication and marketing strategies.

Community Visualization (Girvan-Newman)

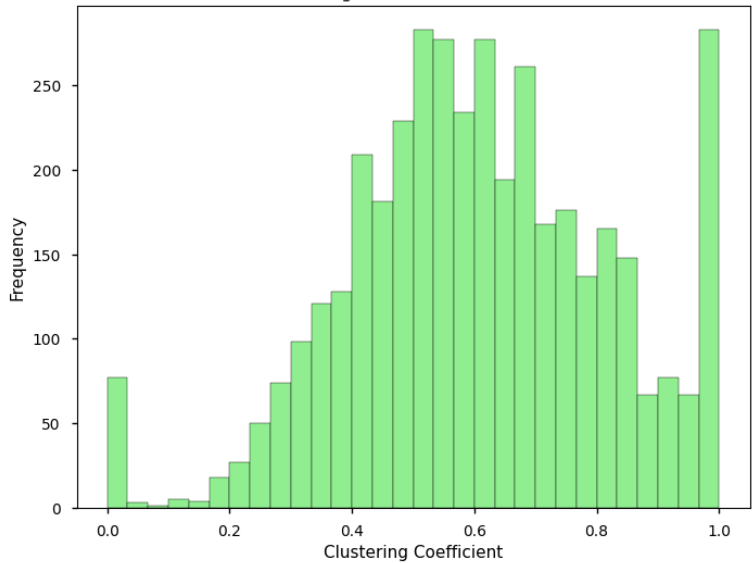


**Clustering Coefficient and K-Core Decomposition**

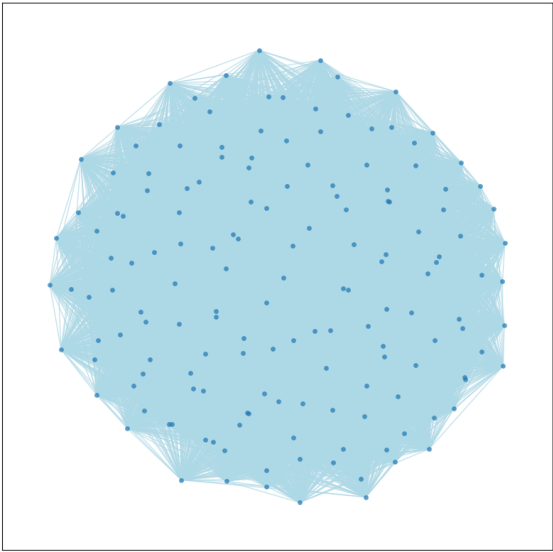
The global clustering coefficient of 0.605 indicates robust tendency for nodes to cluster together, suggesting users form tight-knit groups or communities. High clustering coefficients implies users are more likely to share friends, a characteristic of social networks where people often connect with those who are already connected to them. This behavior enhances the speed and efficiency of information dissemination within these subgroups.

The k-core decomposition analysis, particularly for  $k=1$ , retained nodes with at least one connection, highlighting a resilient network structure. A substantial portion of the network consists of strongly connected subgroups, suggesting resilient community structures. Identification of these groups suggests that they can withstand removal of less connected nodes without significantly compromising overall connectivity of the network.

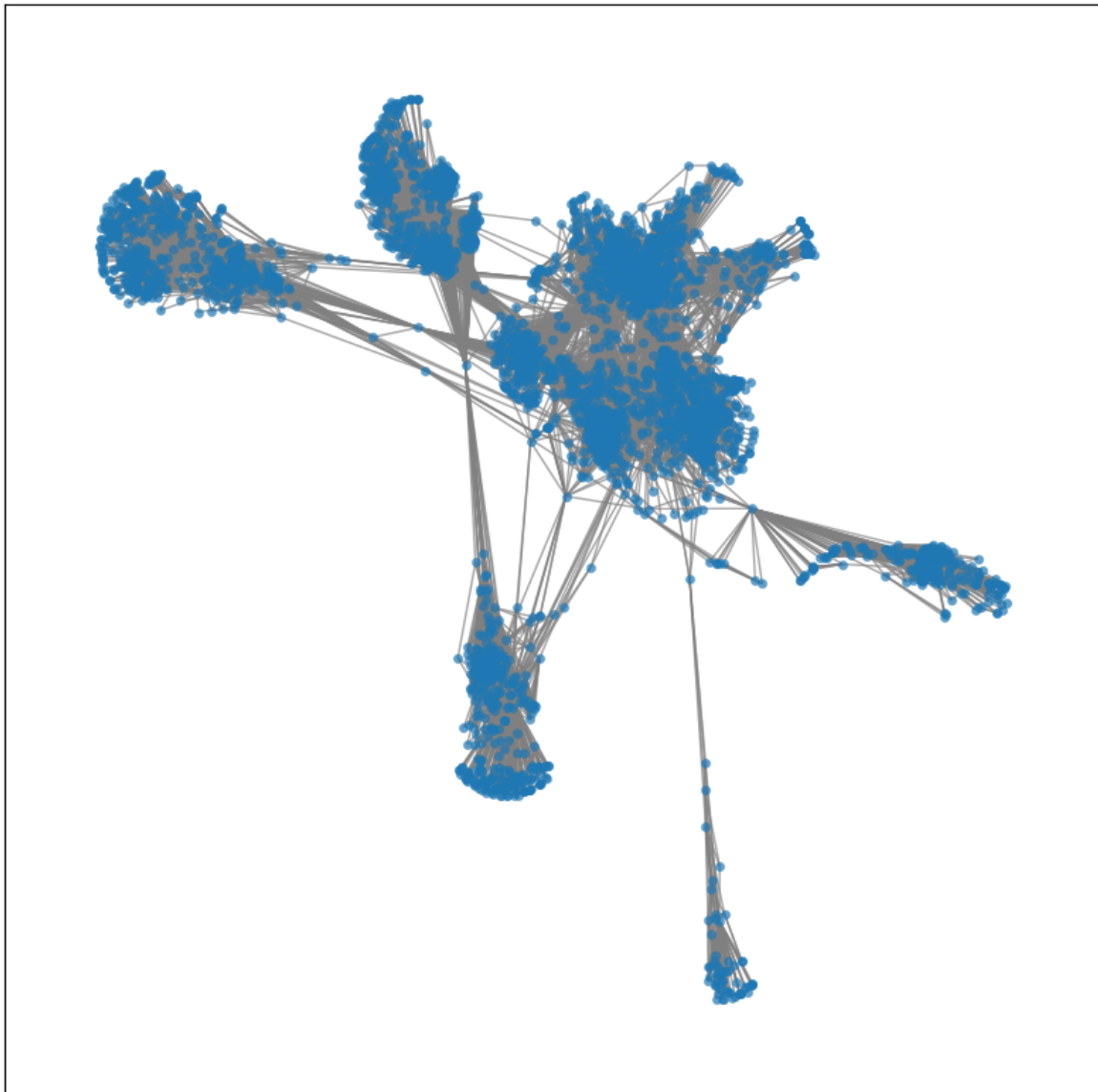
Clustering Coefficient Distribution



k-Core Subgraph Visualization



Node-Link Diagram



### ***Graph Traversal Techniques***

Techniques, such as Breadth-First Search (BFS) and Dijkstra's Algorithm, facilitated efficient exploration of the network, enabling analysis of shortest paths between nodes.

BFS revealed how information flows through a user's direct friends before reaching broader circles. It emphasizes the network's accessibility and users' quick connections.

Dijkstra's Algorithm offered insights into the shortest paths in terms of weighted connections. While this study primarily focused on unweighted edges, the ability to adapt Dijkstra's for scenarios involving varied path costs suggests potential future applications for studying user interactions where some connections may carry different significance.

The application of these algorithms demonstrated the overall navigability of the Facebook network, supporting the notion that users can efficiently reach one another and share information across the network.

## IV.

## CONCLUSION

Effectively employing network analysis reveals structural patterns and dynamics within the Facebook social network. By applying centrality measures, community detection, and traversal techniques, we gained insights into user interactions, influential nodes, and community structures.

### ***Centrality Measures and Influential Nodes***

The analysis found key influential nodes, or information hubs, within the network. Nodes like 107, which ranked high across several metrics, indicated its importance in enabling rapid information dissemination. These nodes act as bridges between communities, making them valuable targets for applications like targeted marketing and public health campaigns.

### ***Community Detection and Social Circles***

Community detection methods, like Louvain and Label Propagation, successfully identified closely-knit groups. These represent user interactions, creating social circles vital for engagement and collaboration. Targeted communication fosters a sense of belonging and more effective outreach.

### ***Clustering Coefficient and K-Core Decomposition***

Users tend to connect with mutual friends, forming strong, cohesive clusters. The decomposition showed tightly connected subgroups maintain strong connectivity even when peripheral nodes are removed. This aligns with the small-world phenomena, enhancing information flow within clusters.

In conclusion social network analysis is crucial for understanding complex digital platforms. It informs future research and practical applications, like improving communication strategies and studying online community evolution.

## V.

## REFERENCES

1. Freeman, L.C. (1978) - "Centrality in social networks: Conceptual clarification." *Social Networks*, 1(3), 215-239.
  2. Newman, M.E.J. (2003). "The structure and function of complex networks." *SIAM Review*, 45(2), 167-256.
  3. Girvan, M., & Newman, M.E.J. (2002). "Community structure in social and biological networks." *Proceedings of the National Academy of Sciences*, 99(12), 7821-7826.
  4. Bonacich, P. (1987). "Power and centrality: A family of measures." *American Journal of Sociology*, 92(5), 1170-1182.
  5. Clauset, A., Newman, M.E.J., & Moore, C. (2004). "Finding community structure in very large networks." *Physical Review E*, 70(6), 066111.
  6. Page, L., Brin, S., Motwani, R., & Winograd, T. (1999). "The PageRank citation ranking: Bringing order to the web." *Stanford InfoLab Technical Report*.
  7. Watts, D.J., & Strogatz, S.H. (1998). "Collective dynamics of 'small-world' networks." *Nature*, 393, 440-442.
  8. Barabási, A.L., & Albert, R. (1999). "Emergence of scaling in random networks." *Science*, 286(5439), 509-512.
  9. Lambiotte, R., Delvenne, J.-C., & Barahona, M. (2008). "Laplacian dynamics and multiscale modular structure in networks." *IEEE Transactions on Network Science and Engineering*, 1(2), 76-90.
  10. Borgatti, S.P., & Everett, M.G. (2006). "A Graph-Theoretic Perspective on Centrality." *Social Networks*, 28(4), 466-484.
-