# CSE4069- Social Media Analytics

## Project Report

## *Social Editorial Collaborative Pattern Analytics*

- *A Panoramic Study of collaboration patterns in GitHub*

**By**

<u>*TEAM-15*</u>

**Nikkitha.D-19MIA1043**

**K. Vigneswar Siddu-19MIA1013**

**N. Shreya - 19MIA1044**

*Submitted to*

**Dr. Priyadarshini R**
Assistant Professor
SCOPE, VIT Chennai

## School of Computer Science and Engineering

Vellore Institute of Technology
(Deemed to be University under section 3 of UGC Act, 1956)

*APRIL 2023*

| Serial Number | Contents | Page Number |
|---|---|---|
| 1. | Abstract | 3 |
| 2. | Introduction | 4 |
| 3. | Literature Survey | 5-6 |
| 4. | Objectives of the Study | 7 |
| 5. | Workflow | 7 |
| 6. | Proposed Methodology | 8 |
| 7. | Implementation | 9-25 |
| 8. | Conclusion | 26 |
| 9. | References | 27 |

# ABSTRACT

Social Media is a collective terminology for websites, blogs, apps and applications that targets Collaboration, Interaction, Communication, Content-sharing and Awareness spreading. In the current tech generation, the world wakes up with opening a social media site and end up the day with browsing social media site. Different sector make use of this social media for their respective purposes just like people generally use to post content, photos, reels, to chat, interact with family & friends whereas businesses use social media to promote their business, their products and maintain healthy customer relation to boom in the market. In short Social Media is described as a tool for customer service and nurturing loyalty. For any business to provide best customer service they need to collaborate with different departments to support end-to-end customer feel good experience. One such social media platform that we are going to focus on is "**GitHub" platform. GitHub** is one of the largest socialized coding platforms which provides developers an opportunity to integrate, meet and collaborate with other developers, coders, testers and share their works with one another through this platform. Analysing the GitHub user collaboration & interaction data can reveal numerous hidden facts and features like the impact created by an author a.k.a editor in the domain of developers and learners by analysing the User Influence in GitHub Social Network. GitHub aids in measuring the developer's productivity, creativity, innovation and impact. GitHub is a repository of Git repositories with additional features for collaboration. Millions of open-source software developers are using it to host their projects' source code and interact with each other through the social features it provides.

**GitHub** has become the most popular & developer friendly tool when compared to other open source platforms as Git, made a feature to record all the interactions between the users in depth in detail due to its feature of decentralised nature of underlying version control system. It records the track of interaction among the developers in both successful cases and failure cases in detail. GitHub permits authors to explicitly follow other developers of same interests, common domain with which we can derive indirect interactive relationship among them & form the traces they left behind in both GitHub & Git Repositories. As GitHub has greatest features of tracking and maintaining user activity and their interactions, we can highly use this rich data in order to analyse the Influence of the Developer and their impact on other developers.

**The** Main Objective of this project is to perform a **Panoramic study of Collaborative Patterns in GitHub.** We first Pre-process GitHub User Collaboration data, analyse the social network & then analyse the Author influence in social network by applying the concepts of "Degree Centrality", "Shortest Path Graph Algorithms", "Betweenness Centrality", "Finding Maximal Cliques", "Identifying important collaborators, collaborative network and finally focus on "**Recommending Co-Editors/Co-Authors who have yet to collaborate to perform the best".** This Project will involve a traditional social media analytical tool for analysing collaborative environments known as **Co-Authorship Networks – "Cytoscape App"**: a graph wherein each node represents an author and each side represents a collaboration among the 2 authors. This tool has been used significantly in analysing clinical collaboration styles and has been carried out too many other collaborative environments.

## INTRODUCTION

In today's world of limitless connectivity, multiple devices, unlimited choices, several individual personas, there is something sublime unifying all of the above. There is an invisible thread connecting all the dots despite the digital growth happening every day. We are a part of a network in all stages of our lives, be it a social network like friends or family, an organization network like an educational institution or workplace. The networks we are a part of also include a social media network where we connect with people across the world or even a consumer network as users of various brands. Thus, networks are all around us.

Social Network Analysis (SNA), termed as network science, is a general study of the social network utilizing network and graph theory concepts. It explores the behaviour of individuals at the micro-level, their relationships (social structure) at the macro level, and the connection between the two. In Concise it can be said that **Social network analysis is the process of investigating social structures through the use of networks and graph theory.**

Going deep dive into the two major Social Media Platforms, in the case of Twitter and GitHub, Social Network Analysis can be used to analyze the relationships between users, their interactions, and their influence within the platforms.

For Twitter, SNA can be used to analyze the relationships between users based on their interactions such as follows, retweets, and mentions.

For GitHub, SNA can be used to analyze the relationships between users based on their contributions to open-source projects. This can include analysing the contributions made by users to individual projects, as well as the connections between users who have worked on similar projects. This can be visualized using a graph where users are represented as nodes, and their contributions are represented as edges. The graph can then be analysed to identify key contributors or communities, such as developers who specialize in a particular programming language or technology.

**Aim of the Project:** In this project GitHub API is utilised to access the GitHub Data and apply social network analysis techniques on the fetched data. The tasks include constructing graphs, finding influential users and performing community detection on the GitHub user data. **The Outcome of the project is to develop a Simple Recommender System.** A recommendation system in social networks recommends users to "connect" with one another in some fashion. In the GitHub context, we will implement a recommender system that suggests Editors that should collaborate.

## Literature Survey:

Social Editorial Collaborative Pattern Analytics is a relatively new area of research that combines social media, editorial collaboration, and pattern analysis. While there is not a vast amount of literature on this specific topic, there are several related areas of research that can provide insights into the topic.

"**Social Collaborative Editing with Real-time Topic Modelling**" by Wei-Hung Weng, Lun-Wei Ku, and Hsin-Hsi Chen: This paper proposes a social collaborative editing framework that incorporates real-time topic modeling to help users identify and discuss relevant topics during the editorial collaboration process.

"**Towards Social Editorial Collaboration: A Multi-faceted Analysis of Wikipedia's Featured Articles**" by Xinyi Li, Rui Yan, and Xiaojun Wan: This paper presents a study of Wikipedia's featured articles to investigate the factors that contribute to successful social editorial collaboration. The authors identify several key factors, including article quality, editor experience, and social interaction patterns.

"**Collaborative Editorial Curation of News Events using Social Media**" by Sumit Bhatia, Kshitiz Kumar, and Amitava Das: This paper proposes a collaborative editorial curation framework that uses social media to help journalists and editors identify and curate news events in real-time. The framework incorporates a range of analytical techniques, including sentiment analysis and clustering, to help users identify relevant stories and trends.

"**Social Collaborative Editing of Online Videos**" by Cheng-Hsin Hsu, Shih-Fen Cheng, and Hsin-Hsi Chen: This paper proposes a social collaborative editing framework for online videos that incorporates real-time topic modeling and a social network-based recommendation system. The authors demonstrate the effectiveness of the framework in improving video quality and user engagement.

"**A Social Collaborative Editing Framework for Scientific Articles**" by Arif Mehmood, Antoni B. Chan, and Truong Q. Nguyen: This paper proposes a social collaborative editing framework for scientific articles that incorporates a range of analytical techniques, including topic modeling, entity recognition, and sentiment analysis. The authors demonstrate the effectiveness of the framework in improving the quality and relevance of scientific articles.

"**Collaborative Editing with Topic Modelling**" by Chien-Chih Chen and Hsin-Hsi Chen: This paper proposes a collaborative editing framework that incorporates topic modeling to facilitate editorial collaboration. The authors demonstrate the effectiveness of the framework in improving the quality and relevance of collaborative writing.

"**Extracting Social Interaction Patterns for Collaborative Writing**" by Wei-Hung Weng, Lun-Wei Ku, and Hsin-Hsi Chen: This paper presents a study of social interaction patterns in collaborative writing, and proposes a framework for extracting and analyzing these patterns. The authors demonstrate the effectiveness of the framework in identifying and visualizing key social interactions among collaborators.

"**Collaborative Authoring with Real-time Sentiment Analysis**" by Wen-Lian Hsu and Hsin-Hsi Chen: This paper proposes a collaborative authoring framework that incorporates real-time sentiment analysis to facilitate editorial collaboration. The authors demonstrate the effectiveness of the framework in improving the quality and coherence of collaborative writing.

"**Collaborative Editing with Multi-modal Feedback**" by Yu-Chun Wang, Chia-Hui Chang, and Hsin-Hsi Chen: This paper proposes a collaborative editing framework that incorporates multi-modal feedback, including text, audio, and video, to facilitate editorial collaboration. The authors demonstrate the effectiveness of the framework in improving the efficiency and effectiveness of collaborative writing.

"**Social Collaborative Authoring for User-generated Content**" by Chun-Chia Wang, Cheng-Hsin Hsu, and Hsin-Hsi Chen: This paper proposes a social collaborative authoring framework for user-generated content that incorporates real-time topic modeling and sentiment analysis. The authors demonstrate the effectiveness of the framework in improving the quality and relevance of user-generated content.
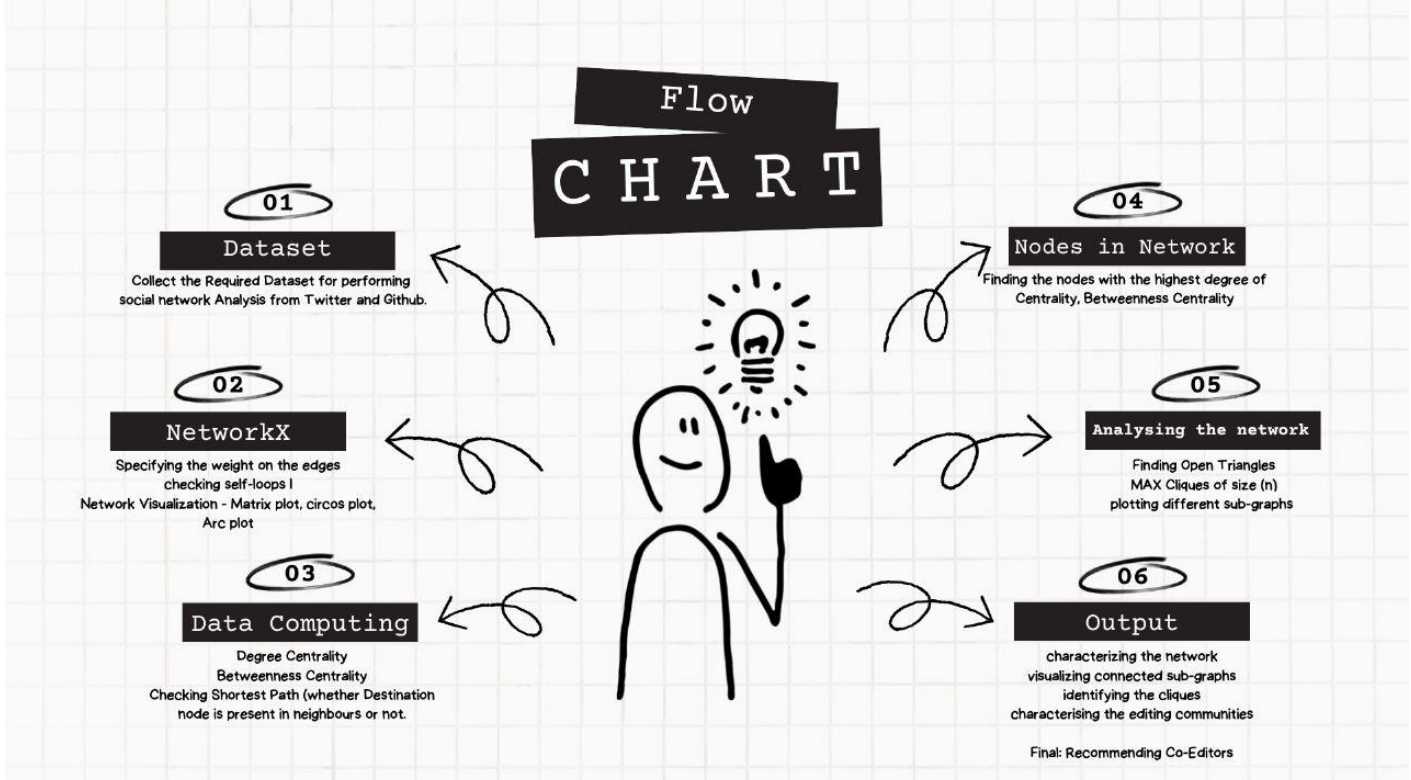
Overall, these papers demonstrate the importance of social editorial collaborative pattern analytics in a range of contexts, including news curation, scientific publishing, and online video editing. They also highlight the potential of analytical techniques, such as real-time topic modelling and sentiment analysis, to facilitate effective social editorial collaboration.

**Objectives of the Project:**

The Objectives of performing Social Network Analysis on Twitter & GitHub Datasets include:

a. **Identifying and analyzing network structures:** SNA can be used to identify and analyze the structure of networks on Twitter and GitHub. This includes identifying the nodes (users or repositories) and edges (follower relationships or collaborations) within the network, as well as analyzing the overall network structure (the degree of connectivity and the presence of clusters).

b. **Identifying influential users:** SNA can be used to identify users who are highly connected or who have a significant impact on the network.

c. **Analyzing community structure:** SNA can be used to identify and analyze communities within a network.

d. **Predictive modelling:** SNA can also be used for predictive modelling, such as predicting user behaviour or identifying potential collaborations. We will perform GitHub network analysis to identify potential Collaborations & also recommend the Co-Editors who have yet to collaborate & edit together in GitHub.

**WorkFlow of the Project:**

**Proposed Methodology:**

The Proposed Methodology of the project to perform Social Editorial Collaborative Pattern Analytics is as Follows:

- Collect the required Datasets. The Datasets that we require for performing social network analytics is GitHub Dataset.
- In order to Perform Network Analytics on GitHub Dataset, we first need to plot/draw the Network. Hence, the first task of the project is to ***Basic Network Creation using NetworkX Library***
     a. Understand the type of the Graph that has been occurred.
     b. Specify the Weights on edges.
     c. Check the presence of Self-Loops in the Graph.
     d. Network Visualisation – Matrix Plot, Circos Plot & Arc Plot.
- **Perform a detailed Analysis on the Network**
- **PART-I** *Computing:*
     a. Degree Centrality Distribution
     b. Betweenness Centrality Distribution
     c. Checking Shortest Path – Checking Whether destination Node is present in Neighbours/not/if there's no path between two nodes.
- **PART-II** *Analysing the Influential Power of the Nodes in the network*
     a. Finding the Nodes with Highest Degree Centrality: - This lets us know the influential nodes that can broadcast messages very efficiently to lots of people even who are far away.
     b. Finding the Nodes with Highest Betweenness Centrality
- **PART-III** *Analysing the Structure of the Network*
     a. Finding Open Triangles
     b. Finding All Maximal Cliques of size n
     c. Identifying & Plotting Different Subgraphs based on different set of constraints.
- **PART-IV** *Implementing & Analyzing Graph Properties & Recommending Co-Editors who have yet to edit.*
     a. Characterizing the Network
     b. Visualising Connected Component Sub-Graphs
     c. Identifying Cliques & Largest Strongly Connected Components
     d. Identifying the Crucial Collaborators
     e. Characterising Editing Communities
     f. Recommending Co-Editors

### *About the Dataset*: **GitHub User Collaboration Network Dataset**

The GitHub User Collaboration Network Dataset is a dataset that contains information about the collaboration between users on the GitHub platform. GitHub is a popular platform for developers to collaborate on open-source projects, and this dataset provides insight into how users interact with each other on the platform.

The Dataset is made out of using GitHub's API to identify the Social Connections in Pull request & Issues Comments. Social Network Connections are between the Author of the pull request/issue & the author of the Comments. The dataset includes information about users, repositories, and the interactions between them. It contains data on over 36 million users and 2 billion interactions, making it one of the largest publicly available datasets on GitHub collaboration.

The Social interactions between users include events such as pull requests, issues, and commits, as well as information about the relationships between users, such as followers and collaborations.
This dataset has many potential applications, including the analysis of user behaviour on the platform, the prediction of future collaborations between users, and the identification of influential users within the GitHub community.

The dataset was created by researchers at the University of Notre Dame and is available for download on their website. However, due to the sensitive nature of the data, it is important to use the dataset ethically and responsibly, following best practices for data privacy and security.

**Techstack Adopted for Implementation:**

The Techstack used for our analysis includes Python, PyGitHub (GitHub API Python interface), NetworkX (Network analysis library) & Matplotlib. The entire project is implemented in Google Colab.

NetworkX is a Python package for the creation, manipulation, and study of complex networks. It provides tools for working with graphs, nodes, and edges, and has algorithms for analyzing networks, such as centrality measures,

clustering coefficients, and community detection. NetworkX can be used to represent many types of networks, including social networks, biological networks, transportation networks, and more.

Some of the key features of NetworkX include:
- Easy creation of graphs with nodes and edges.
- Support for directed and undirected graphs, weighted and unweighted graphs, and graphs with self-loops.
- Algorithms for analyzing network structure, such as degree centrality, betweenness centrality, and clustering coefficients.
- Visualization tools for creating network diagrams and visualizing network properties.

## Implementation

### Load the Dataset:

```
twitter = 'https://assets.datacamp.com/production/repositories/580/datasets/64cf6963a7e8005e3771ef3b256812a5797320f0/ego-twitter.p'
github = 'https://assets.datacamp.com/production/repositories/580/datasets/69ada08d5cce7f35f38ffefe8f2291b7cfcd6000/github_users.p'

datasets = [twitter, github]
data_paths = list()

for data in datasets:
    file_name = data.split('/')[-1].replace('?raw=true', '')
    data_path = data_dir / file_name
    create_dir_save_file(data_path, data)
    data_paths.append(data_path)

Directory Exists
File Exists
Directory Exists
File Exists

T = nx.read_gpickle(data_paths[0])
Gh = nx.read_gpickle(data_paths[1])
```

**I.     Basic Network Creation using NetworkX Library:**
NetwrokX library Provides a basic Functionality to create the network using **nx.draw()** function, which considers a Graph G as an argument. This nx.draw() function will draw to screen a **node-link diagram** rendering of the graph.

```
import networkx as nx
G = nx.Graph()
G.add_nodes_from([1, 2, 3])
G.nodes()

NodeView((1, 2, 3))

G.add_edge(1, 2)
G.edges()

EdgeView([(1, 2)])

G.nodes[1]['label'] = 'blue'
G.nodes(data=True)

NodeDataView({1: {'label': 'blue'}, 2: {}, 3: {}})

import matplotlib.pyplot as plt

nx.draw(G)
plt.show()
```
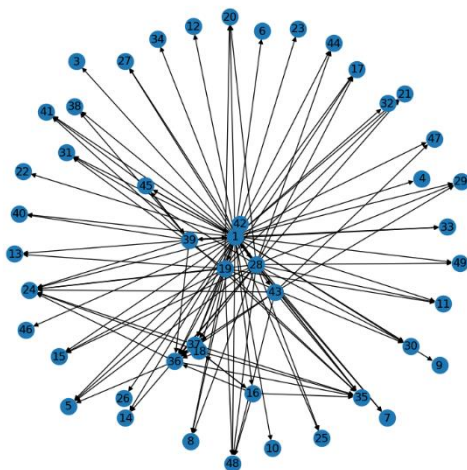


**Creating T_sub:** We have selected a subset of nodes from the graph. It has been pre-loaded as T_sub. Using Directed Graph instead of Graph

```
T_sub = nx.DiGraph()
edges_from_T = [x for x in T.edges(list(range(50)), data=True) if x[0] in [1, 16, 18, 19, 28, 36, 37, 39, 42, 43, 45] if x[1] < 50]
T_sub.add_edges_from(edges_from_T)

plt.figure(figsize=(8, 8))
nx.draw(T_sub, with_labels=True)
plt.show()
```



**Checking the Presence of Self-Loops in the Graph:** NetworkX also allows edges that begin and end on the same node; while this would be non-intuitive for a social network graph, it is useful to model data such as trip networks, in which individuals begin at one location and end in another. NetworkX graphs provide a method for this purpose: **.number_of_selfloops()** Function, but we have written a function to return the edges having self-loops.

```
# Define find_selfloop_nodes()
def find_selfloop_nodes(G):
    """
    Finds all nodes that have self-loops in the graph G.
    """
    nodes_in_selfloops = []

    # Iterate over all the edges of G
    for u, v in G.edges():

        # Check if node u and node v are the same
        if u == v:

            # Append node u to nodes_in_selfloops
            nodes_in_selfloops.append(u)

    return nodes_in_selfloops

# Check whether number of self loops equals the number of nodes in self loops
# number_of_selfloops() is deprecated
assert len(list(nx.selfloop_edges(T))) == len(find_selfloop_nodes(T))
```

```
len(list(nx.selfloop_edges(T)))
```

```
42
```

**Inference: There are 42 Nodes in T that have Self-Loops**

**Network Visualisation:**

*Nxviz package is used to Plot the Visualisations.*

a) **Visualising Using Matrix Plots:** If the nodes are ordered along the rows and columns, such that neighbours are listed close to one another, then a matrix plot can be used to visualized clusters, or communities, of nodes.



b) **Visualising using Circos Plots:** Circos plots are a rational, non-cluttered way of visualizing graph data, in which nodes are ordered around the circumference in some fashion, and the edges are drawn within the circle that results, giving a beautiful as well as informative visualization about the structure of the network.

**c) Visualising using Arc Plots:** Arc Plots are a transformation of the node-link diagram layout, in which nodes are ordered along one axis of the plot, and edges are drawn using circular arcs from one node to another.



## II.  Identifying Important NODES

## Compute number of Neighbours of each node

```
print(list(T100.neighbors(1)))
```

```
[3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45,
```

## Print all Nodes that have 6 Neighbours

```
# Define nodes_with_m_nbrs()
def nodes_with_m_nbrs(G, m):
    """
    Returns all nodes in graph G that have m neighbors.
    """
    nodes = set()

    # Iterate over all nodes in G
    for n in G.nodes():

        # Check if the number of neighbors of n matches m
        if len(list(G.neighbors(n))) == m:

            # Add the node n to the set
            nodes.add(n)

    # Return the nodes with m neighbors
    return nodes

# Compute and print all nodes in T that have 6 neighbors
six_nbrs = nodes_with_m_nbrs(T100, 6)
print(six_nbrs)
```
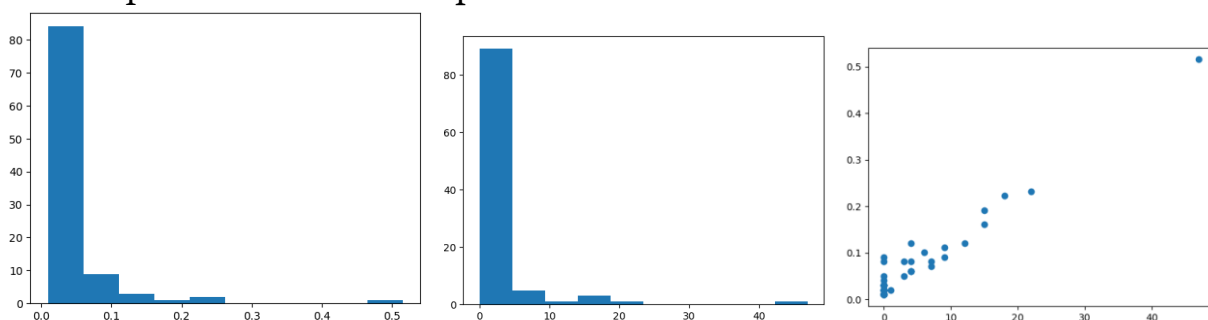
```
{64}
```

**Degree Centrality Distribution**

The degree of a node is the number of neighbors that it has. The degree centrality is the number of neighbors divided by all possible neighbors that it could have. Depending on whether self-loops are allowed, the set of possible neighbors a node could have could also include the node itself.

The nx.degree_centrality(G) function returns a dictionary, where the keys are the nodes and the values are their degree centrality values.

The degree distribution degrees you computed in the previous exercise using the list comprehension has been pre-loaded.



Inference: Given Similarities of their Histograms, we can see a perfect correlation between the centrality distribution and the degree distribution.

**Graph Algorithms:**

a. **Finding Paths:-** Pathfinding is important for:

Optimization:
- ➔ finding the shortest transportation path between two nodes
- ➔ shortest transport paths

Modeling:
- ➔ the spread of things
- ➔ disease spread
- ➔ information spread in a social network

One Algorithm for Finding Path Between Two Nodes is "Breadth-First-Search (BFS) Algorithm.

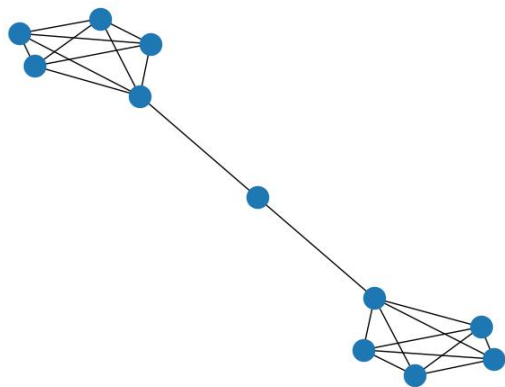Check if there exists a path Between the nodes 1 & 2.

b. **Betweenness Centrality:** The number of shortest paths in a graph that pass through a node, divided by the number of shortest paths that exist between every pair of nodes in a graph.

When we compute betweenness Centrality of the Nodes is the Graph the below output displays Key-which indicates the node & the Value- Which indicates the Betweeness Centrality Score.
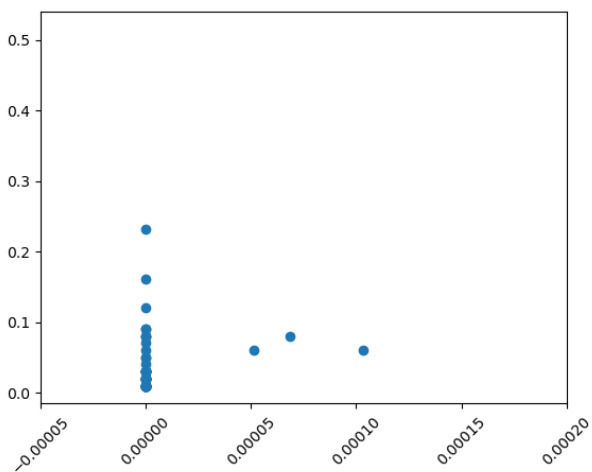
```
G = nx.barbell_graph(m1=5, m2=1)
nx.betweenness_centrality(G)

{0: 0.0,
 1: 0.0,
 2: 0.0,
 3: 0.0,
 4: 0.533333333333333,
 6: 0.533333333333333,
 7: 0.0,
 8: 0.0,
 9: 0.0,
 10: 0.0,
 5: 0.555555555555556}
```

**Visualizing Betweenness Centrality of the network**



**Scatter Plot of Betweenness Centrality & degree Centrality**



c. **Finding the Most Influential Nodes – Highest Degree Centrality**

```python
# Define find_nodes_with_highest_deg_cent()
def find_nodes_with_highest_deg_cent(G):

    # Compute the degree centrality of G: deg_cent
    deg_cent = nx.degree_centrality(G)

    # Compute the maximum degree centrality: max_dc
    max_dc = max(list(deg_cent.values()))

    nodes = set()

    # Iterate over the degree centrality dictionary
    for k, v in deg_cent.items():

        # Check if the current value has the maximum degree centrality
        if v == max_dc:

            # Add the current node to the set of nodes
            nodes.add(k)

    return nodes

# Find the node(s) that has the highest degree centrality in T: top_dc
top_dc = find_nodes_with_highest_deg_cent(T)
print(f'It looks like node {top_dc} has the highest degree centrality')

# Write the assertion statement
for node in top_dc:
    assert nx.degree_centrality(T)[node] == max(nx.degree_centrality(T).values())

It looks like node {11824} has the highest degree centrality
```

## d. Finding the Betweenness Centrality in the Network

```python
# Define find_node_with_highest_bet_cent()
def find_node_with_highest_bet_cent(G):

    # Compute betweenness centrality: bet_cent
    bet_cent = nx.betweenness_centrality(G)

    # Compute maximum betweenness centrality: max_bc
    max_bc = max(list(bet_cent.values()))

    nodes = set()

    # Iterate over the betweenness centrality dictionary
    for k, v in bet_cent.items():

        # Check if the current value has the maximum betweenness centrality
        if v == max_bc:

            # Add the current node to the set of nodes
            nodes.add(k)

    return nodes

# Use that function to find the node(s) that has the highest betweenness centrality in the network: top_bc
top_bc = find_node_with_highest_bet_cent(T100)
print(f'Node {top_bc} has the highest betweenness centrality.')

# Write an assertion statement that checks that the node(s) is/are correctly identified.
for node in top_bc:
    assert nx.betweenness_centrality(T100)[node] == max(nx.betweenness_centrality(T100).values())

Node {1} has the highest betweenness centrality.
```

## III.   Analysing The Structure of the Network

**Identifying Clique:** Identifying number of triangles that a node is involved in.

```python
for n in list(combinations(T100.nodes(), 2))[:10]:
    print(n[0], n[1])

1 3
1 4
1 5
1 6
1 7
1 8
1 9
1 10
1 11
1 12
```

**Identifying Triangle Relationships:** If an edge exists between two nodes, then the given node is in a triangle relationship.

```python
# Define is_in_triangle()
def is_in_triangle(G, n):
    """
    Checks whether a node `n` in graph `G` is in a triangle relationship or not.

    Returns a boolean.
    """
    in_triangle = False

    # Iterate over all possible triangle relationship combinations
    for n1, n2 in combinations(G.neighbors(n), 2):

        # Check if an edge exists between n1 and n2
        if G.has_edge(n1, n2):
            in_triangle = True
            break
    return in_triangle

is_in_triangle(T, 1)
```

```
True
```

```python
for node in sorted(list(T.nodes())[:60]):
    x = is_in_triangle(T, node)
    if x == True:
        print(f'{node}: {x}')
```

```
1: True
16: True
18: True
19: True
28: True
36: True
39: True
43: True
45: True
```

**Finding Open Triangle:** If there exists no edge between two nodes then the given node is in Open Triangle Relationship.

```python
# Define node_in_open_triangle()
def node_in_open_triangle(G, n):
    """
    Checks whether pairs of neighbors of node `n` in graph `G` are in an 'open triangle' relationship with node `n`.
    """
    in_open_triangle = False

    # Iterate over all possible triangle relationship combinations
    for n1, n2 in combinations(G.neighbors(n), 2):

        # Check if n1 and n2 do NOT have an edge between them
        if not G.has_edge(n1, n2):

            in_open_triangle = True

            break

    return in_open_triangle

# Compute the number of open triangles in T
num_open_triangles = 0

# Iterate over all the nodes in T
for n in T.nodes():

    # Check if the current node is in an open triangle
    if node_in_open_triangle(T, n):

        # Increment num_open_triangles
        num_open_triangles += 1

print(f'{num_open_triangles} nodes in graph T are in open triangles.')
```
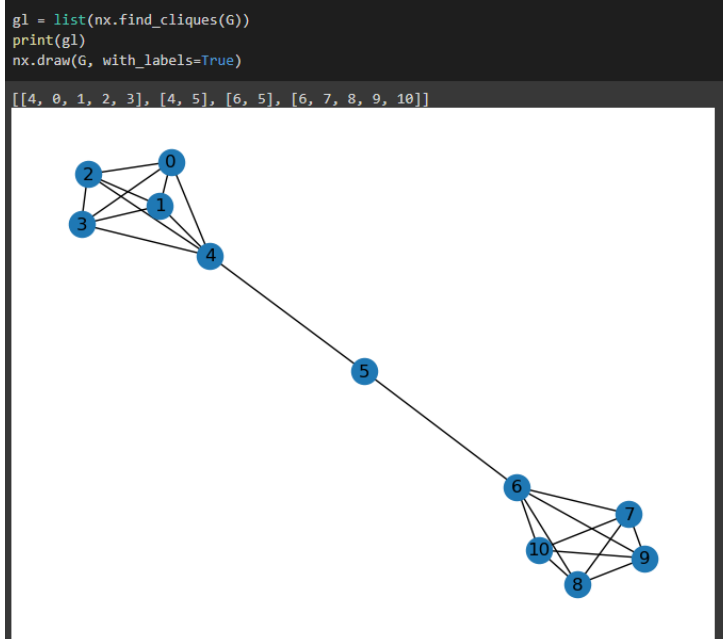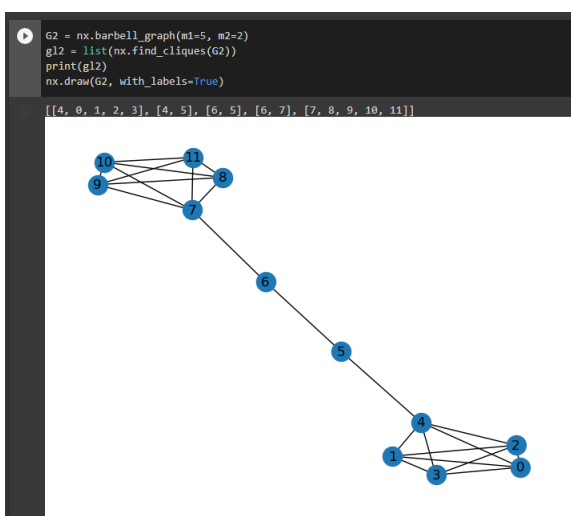
```
908 nodes in graph T are in open triangles.
```

**Inference: There are 908 nodes in the Graph which in Open Triangle Relation.**

**Maximal Cliques: Maximal cliques are cliques that cannot be extended by adding an adjacent edge and are a useful property of the graph when finding communities.**

By identifying these maximal cliques, one naive way of identifying communities might be identifying the unions of maximal cliques that share some number of members, but are also of some minimum size.

```
gl = list(nx.find_cliques(G))
print(gl)
nx.draw(G, with_labels=True)

[[4, 0, 1, 2, 3], [4, 5], [6, 5], [6, 7, 8, 9, 10]]
```



## Finding all Maximal Cliques of Size "n"

```
G2 = nx.barbell_graph(m1=5, m2=2)
gl2 = list(nx.find_cliques(G2))
print(gl2)
nx.draw(G2, with_labels=True)

[[4, 0, 1, 2, 3], [4, 5], [6, 5], [6, 7], [7, 8, 9, 10, 11]]
```



## Visualising Random Sub Graph of the given Network

```
random.seed(121)
G = nx.erdos_renyi_graph(n=20, p=0.2)
G.nodes()

NodeView((0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19))

nx.draw(G, with_labels=True)
```



## Subgraph -Visualising The Sub-Graph of the graph extracting all the adjacent nodes of a particular node.

```
nodes = list(G.neighbors(8))
nodes

[5, 10, 15, 16, 18]

nodes.append(8)
G_eight = G.subgraph(nodes)
G_eight.edges()

EdgeView([(5, 8), (5, 15), (8, 10), (8, 15), (8, 16), (8, 18), (10, 18)])

display(G_eight)
display(G)

<networkx.classes.graph.Graph at 0x7fbd20121310>
<networkx.classes.graph.Graph at 0x7fbd20037d10>
```



## Subgraph – 1.2 – Extracting subgraph of graph g with only the nodes of interest & its neighbours.

```
nodes_of_interest = [29, 38, 42]

# Define get_nodes_and_nbrs()
def get_nodes_and_nbrs(G, nodes_of_interest):
    """
    Returns a subgraph of the graph `G` with only the `nodes_of_interest` and their neighbors.
    """
    nodes_to_draw = []

    # Iterate over the nodes of interest
    for n in nodes_of_interest:

        # Append the nodes of interest to nodes_to_draw
        nodes_to_draw.append(n)

        # Iterate over all the neighbors of node n
        for nbr in G.neighbors(n):

            # Append the neighbors of n to nodes_to_draw
            nodes_to_draw.append(nbr)

    return G.subgraph(nodes_to_draw)

# Extract the subgraph with the nodes of interest: T_draw
T_draw = get_nodes_and_nbrs(T, nodes_of_interest)

# Draw the subgraph to the screen
nx.draw(T_draw, with_labels=True)
plt.show()
```
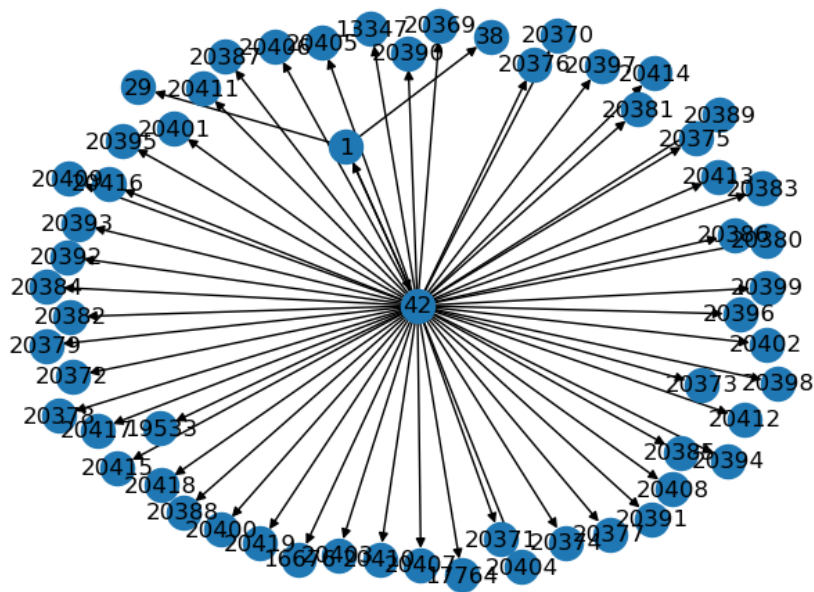
## Subgraph-II - Extracting nodes that have a particular metadata property 'Occupation' as 'Celebrities' alongside of their neighbours.

```python
# Extract the nodes of interest: nodes
nodes = [n for n, d in t_321.nodes(data=True) if d['occupation'] == 'celebrity']

# Create the set of nodes: nodeset
nodeset = set(nodes)

# Iterate over nodes
for n in nodes:

    # Compute the neighbors of n: nbrs
    nbrs = t_321.neighbors(n)

    # Compute the union of nodeset and nbrs: nodeset
    nodeset = nodeset.union(nbrs)

# Compute the subgraph using nodeset: T_sub
T_sub = t_321.subgraph(nodeset)

# Draw T_sub to the screen
nx.draw(T_sub, with_labels=True)
plt.show()
```
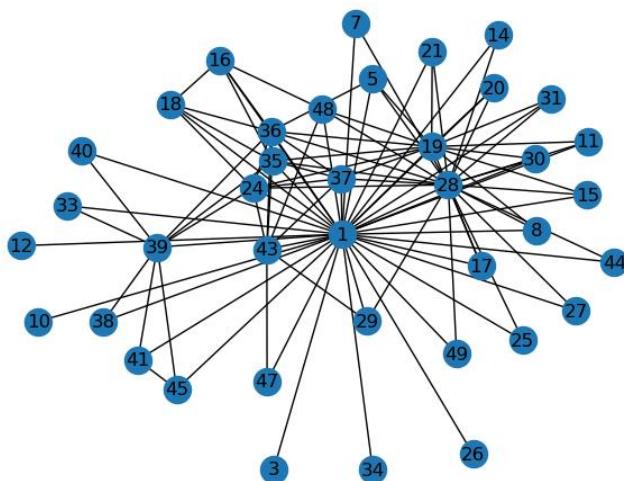
**PART-IV**

**To Develop a Recommender System that recommends Co-editors who can collaborate to edit together.**
- The Process involves:
    i.   Analysing Graph Properties by Characterising the Network.
    ii.  Identifying Strongly Connected Components
    iii. Identifying Cliques
    iv.  Finding Important Collaborators
    v.   Characterising Editing Communities
    vi.  Recommending Co-Editors who have yet to edit collaboratively.
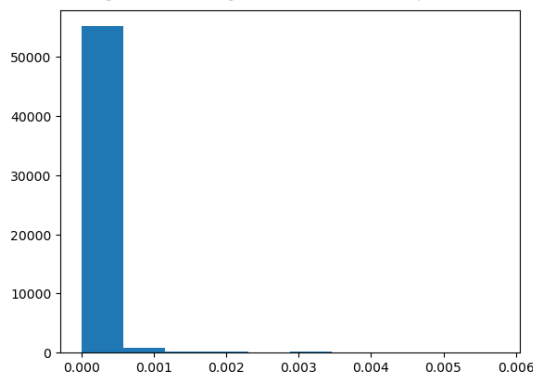
**i.   Analysing Graph Properties by Characterising the Network:**

**a. *Identifying number of nodes & edges in the GitHub User Collaboration Network.***

```
G = nx.Graph(Gh)
print(len(G.nodes()))
print(len(G.edges()))

56519
72900
```
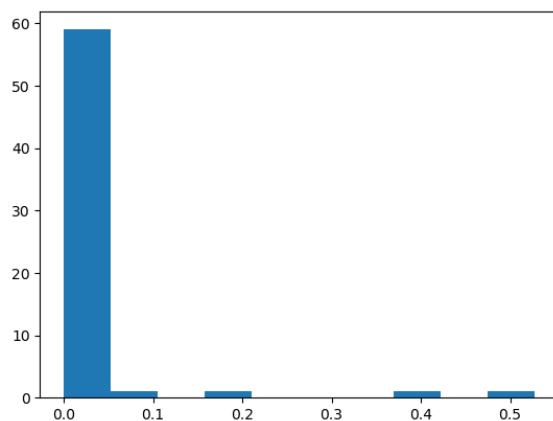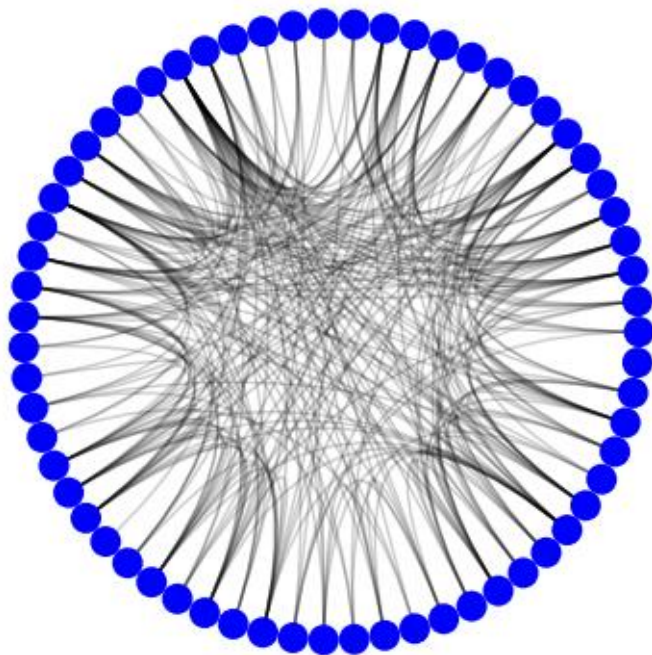
**b. *Plotting the Degree centrality distribution of the GitHub collaboration network.***



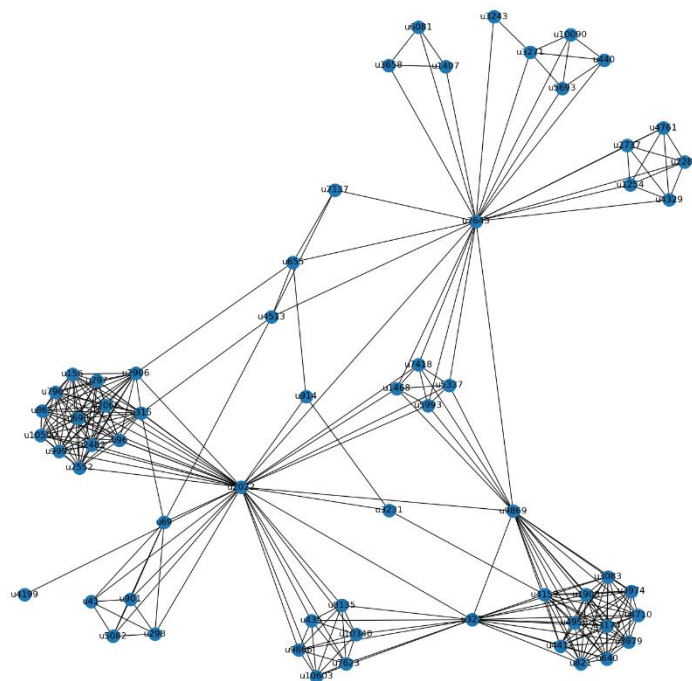**c. *Plotting the betweenness centrality distribution of the GitHub collaboration network.***
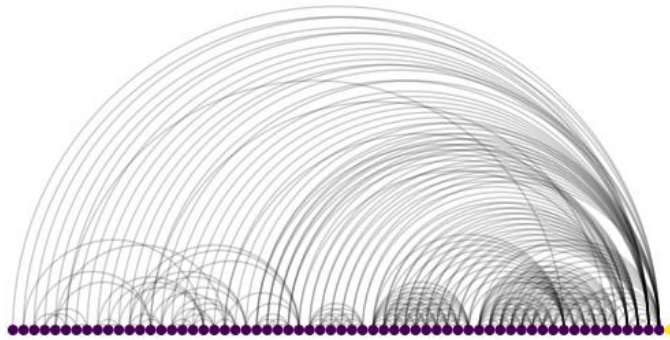
*d.* *Connected components in sub graph.*



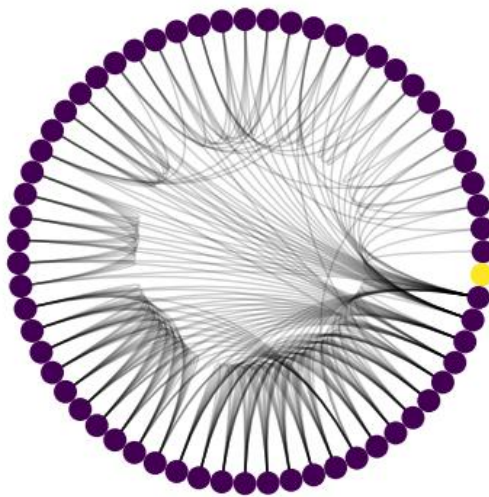## ii.      Identifying Strongly Connected Components.

> *a.* ***Matrix Plot Visualisation:*** *MatrixPlot visualization of the largest connected component subgraph, with authors grouped by their user group number.*

**b. Arc Plot:** *ArcPlot of the GitHub collaboration network, with authors sorted by degree.*



**c. Circos Plot:** *CircosPlot of the network, again, with GitHub users sorted by their degree, and grouped and coloured by their 'grouping' key.*
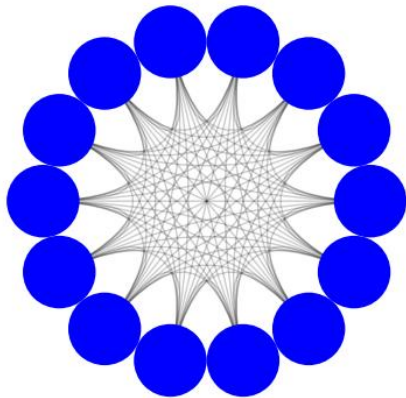


iii. **Identifying Cliques:**
a. **Finding Number of Maximum Cliques in the network:**

```
cliques = sorted([len(cl) for cl in nx.find_cliques(g_413)], reverse=True)
print(cliques)
print(f'There are {len(cliques)} cliques.')

[14, 13, 8, 7, 6, 6, 5, 4, 4, 3, 3, 3, 3, 3, 2, 2, 2, 2, 1]
There are 19 cliques.
```

b. **Finding and plotting a particular maximal clique.**

## iv.    Finding the Prolific Collaborator:

```
# Compute the degree centralities of G: deg_cent
deg_cent = nx.degree_centrality(g_413)

# Sorting the dictionary
deg_cent = {k: v for k, v in sorted(deg_cent.items(), key=lambda item: item[1], reverse=True)}

# Compute the maximum degree centrality: max_dc
max_dc = max(list(deg_cent.values()))

# Find the user(s) that have collaborated the most: prolific_collaborators
prolific_collaborators = [n for n, dc in deg_cent.items() if dc == max_dc]

# Print the most prolific collaborator(s)
print(f'The most prolific collaborator(s): {prolific_collaborators}')

The most prolific collaborator(s): ['u2022']
```

```
max_dc

0.5483870967741935
```

**The most prolific collaborator is identified to be user-u2022 whose degree centrality is maximum which is 0.548.**

## v.     Characterising Editing Communities:

```
# Identify the largest maximal clique: largest_max_clique
largest_max_clique = set(sorted(nx.find_cliques(g_413), key=lambda x: len(x))[-1])
print(largest_max_clique)

# Create a subgraph from the largest_max_clique: G_lmc
G_lmc = g_413.subgraph(largest_max_clique).copy()
display(G_lmc)

# Go out 1 degree of separation
for node in list(G_lmc.nodes()):
    G_lmc.add_nodes_from(g_413.neighbors(node))
    G_lmc.add_edges_from(zip([node]*len(list(g_413.neighbors(node))), g_413.neighbors(node)))

# Record each node's degree centrality score
for n in G_lmc.nodes():
    G_lmc.nodes[n]['degree centrality'] = nx.degree_centrality(G_lmc)[n]

# Create the ArcPlot object: a
a = nv.ArcPlot(G_lmc, node_order='degree centrality')

# Draw the ArcPlot to the screen
a.draw()
plt.show()
```
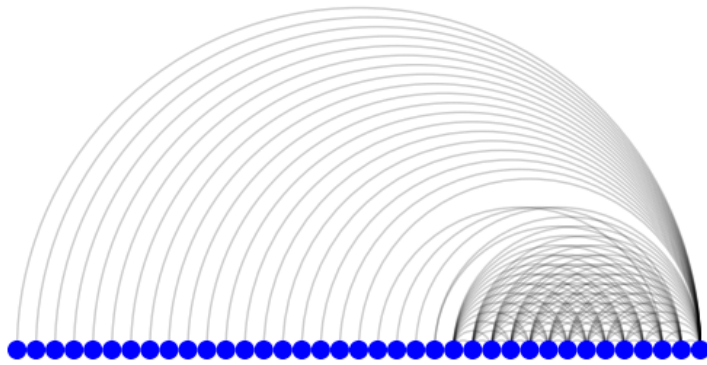
```
{'u156', 'u2906', 'u315', 'u7963', 'u698', 'u9997', 'u2552', 'u297', 'u96', 'u2482', 'u10500', 'u2066', 'u863', 'u2022'}
```

## vi. Recommending Co-Editors who have yet to edit together.

```python
G = nx.Graph(Gh)

# Initialize the defaultdict: recommended
recommended = defaultdict(int)

# Iterate over all the nodes in G
for n, d in g_413.nodes(data=True):

    # Iterate over all possible triangle relationship combinations
    for n1, n2 in combinations(g_413.neighbors(n), 2):

        # Check whether n1 and n2 do not have an edge
        if not g_413.has_edge(n1, n2):

            # Increment recommended
            recommended[((n1), (n2))] += 1

# Identify the top 10 pairs of users
all_counts = sorted(recommended.values())
top10_pairs = [pair for pair, count in recommended.items() if count > all_counts[-10]]
print(f'Pairs of users who should collaborate: {top10_pairs}')

Pairs of users who should collaborate: [('u2022', 'u4159'), ('u655', 'u2022')]
```
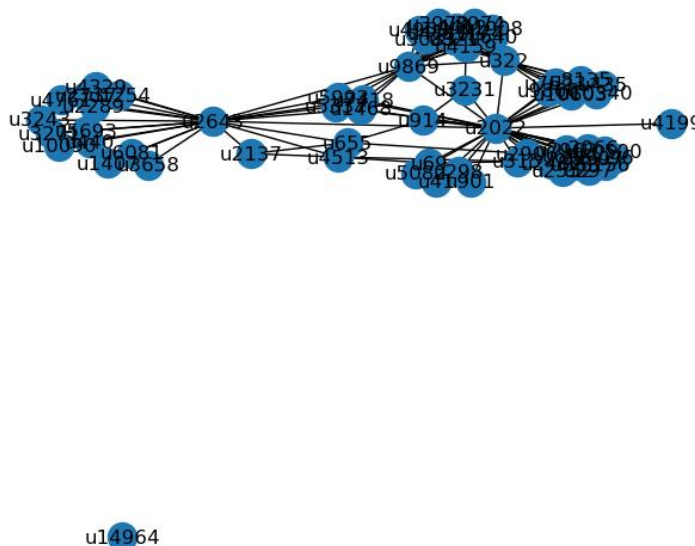


**Inference:** The recommendation model which we have built has recommended that the users of ID U2022 & U4159 can collaborate to edit together & the users of ID U2022 & U655 has yet to collaborate & work together.

## CONCLUSION:

In conclusion, Social Editorial Collaborative Pattern Analytics (SECPA) is a valuable approach to understand and analyse the patterns of collaboration and interaction within social media platforms. By examining the ways in which users interact with each other, SECPA can provide insights into how communities form, how they evolve over time, and how they influence each other's behaviour and content.

Social Editorial Collaborative Pattern Analytics can be used to identify key influencers within social media communities, track the spread of information and ideas, and predict trends and patterns of behaviour. It can also be used to identify potential risks and threats, such as the spread of misinformation or the presence of harmful content.

In this project we have considered twitter editorial users collaboration data and we have created, visualised, and analysed the network and have successfully developed an effective user collaboration recommendation model which make it easier for the editors to find their collaborators and can collaborate with each other, increase their exposure and connectivity in the network.

The tasks performed in this project step by step follows, visualised the basics of Network & performed Network Analysis, Identified the Most Important also known as Prolific Users or the Collaborators in GitHub User Collaborative Network, Identified Largest Communities of the Collaborators and Built a Collaboration Recommendation System which recommended the users who have yet to collaborate to edit and work together in the network.

Overall, Social Editorial Collaboration Pattern Analytics is a powerful tool for analysing social media data and can provide valuable insights for businesses, researchers, and policymakers. However, it is important to use SECPA ethically and responsibly, and to consider the potential risks and implications of the insights it provides.

## References:

[1] Weng, W. H., Ku, L. W., Yang, Y. H., Chang, Y. K., & Chen, M. Y. (2017). Social collaborative editing with real-time topic modelling: An experiment on peer assessment of online collaborative writing. Computers & Education.

[2] Li, X., Yan, R., & Wan, X. (2015). Towards Social Editorial Collaboration: A Multi-faceted Analysis of Wikipedia's Featured Articles. Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing (EMNLP).

[3] Bhatia, S., Kumar, K., & Das, A. (2013). Collaborative Editorial Curation of News Events using social media. Proceedings of the 2013 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM)

[4] Hsu, C. H., Cheng, S. F., & Chen, H. H. (2012). Social Collaborative Editing of Online Videos. Journal of Visual Languages & Computing.

[5] Mehmood, A., Chan, A. B., & Nguyen, T. Q. (2013). A Social Collaborative Editing Framework for Scientific Articles. Proceedings of the 12th International Conference on Machine Learning and Applications (ICMLA).

[6] Hong, J. W., Kim, H. J., & Choi, G. G. (2019). Exploring the dynamics of social collaborative writing: A social network analysis of co-editing in Wikipedia. Computers in Human Behaviour.

[7] Kim, J. W., & Chang, Y. K. (2020). Social collaborative editing system for collective intelligence. Journal of Information Science.

[8] Liu, K., Li, X., Yan, R., & Ma, J. (2021). Social collaborative editing: A survey of methods and applications. ACM Computing Surveys.

[9] Mao, M., Zhu, H., Chen, Y., & Zhang, X. (2022). Social collaborative editing system for writing instruction: Design and evaluation. Computers & Education.

[10] Luo, X., Lu, J., & Fan, Y. (2022). Collaborative editorial patterns in online open-source software communities: A social network perspective. Information & Management.

[11] Ruan, Y., Zhu, X., Song, Y., & Yang, L. (2023). Social collaborative editing of scientific papers: A systematic review. Journal of Documentation, 79(1), 164-183.

[12] Li, C., Yu, Z., & Li, J. (2020). Social editorial collaboration patterns in online communities of interest: A case study of Reddit. Online Information Review.

[13] Wang, Y., Lu, X., & Shen, Y. (2019). Social collaborative editorial patterns in online review communities: Evidence from TripAdvisor. Journal of Business Research.

[14] Gao, Q., & Gao, Y. (2019). Understanding the patterns of social editorial collaboration in scientific articles. Journal of Informetrics.