# Graph Data Mining-based Community Clustering for Gender-Biased Community Detection of Social Media

by

Sanjoy Dev
19101507
Maliha Tabassum
19101212
Mohammad Rahat Khan
20101616
Maliha Bushra Hoque
19101543
Basharat Fatema
20101600

A thesis submitted to the Department of Computer Science and Engineering
in partial fulfillment of the requirements for the degree of
B.Sc. in Computer Science

Department of Computer Science and Engineering
School of Data and Sciences
Brac University
May 2023

# Declaration

It is hereby declared that

1. The thesis submitted is my/our own original work while completing degree at Brac University.

2. The thesis does not contain material previously published or written by a third party, except where this is appropriately cited through full and accurate referencing.

3. The thesis does not contain material which has been accepted, or submitted, for any other degree or diploma at a university or other institution.

4. We have acknowledged all main sources of help.

**Student's Full Name & Signature:**

_Sanjoy Dev_

---
Sanjoy Dev
19101507

_Maliha Tabassum_

---
Maliha Tabassum
19101212

_Mohammad Rahat Khan._

---
Mohammad Rahat Khan
20101616

_Maliha Bushra Hoque_

---
Maliha Bushra Hoque
19101543

_Basharat Fatema_

---
Basharat Fatema
20101600

# Approval

The thesis/project titled "Graph Data Mining-based Community Clustering for Gender-Biased Community Detection of Social Media" submitted by

1. Sanjoy Dev (19101507 )

2. Maliha Tabassum (19101212)

3. Mohammad Rahat Khan (20101616)

4. Maliha Bushra Hoque (19101543)

5. Basharat Fatema (20101600)

Of Spring, 2023 has been accepted as satisfactory in partial fulfillment of the requirement for the degree of B.Sc. in Computer Science on May 22, 2023.

**Examining Committee:**

Supervisor:
(Member)

_____
Md. Golam Rabiul Alam, Ph.D.
Professor
Department of Computer Science and Engineering
BRAC University

Program Coordinator:
(Member)

_____
Md. Golam Rabiul Alam, Ph.D.

Professor
Department of Computer Science and Engineering
BRAC University

Head of Department:
(Chair)

_____
Sadia Hamid Kazi

Associate Professor
Department of Computer Science and Engineering
BRAC University

# Abstract

One of the alarming and uprising issues of the world is gender inequality in recent decades. It is a widespread problem that affects people all around the world, albeit its manifestations and severity vary depending on society and culture. The research gives a thorough investigation of gender bias in online social networks utilizing community clustering and graph data mining approaches. The research methodology includes gathering Twitter data about gender bias using some specific keywords and utilizing networkX to build a graph representation. To divide the graph into different communities, three well-known community detection algorithms—Louvain, Girvan-Newman, and Walktrap—are used. These algorithms' effectiveness is assessed using extrinsic metrics like V-measure and normalized mutual information (NMI), as well as intrinsic metrics like F1 score, recall, and precision. The characteristics of the selected communities are also studied using descriptive statistics and visualization methods. Four communities on gender biasness: Male Biased, Female Biased, Feminism and Neutral people are presented here. The research advances knowledge of gender biases in online social networks and can guide initiatives to advance equality and inclusivity. The goal of this study is to create a solid framework for identifying and examining communities that show neutrality, feminism, neutrality, and male and female prejudice.

# Acknowledgement

Firstly, all praise to the Great Allah for whom our thesis have been completed without any major interruption.

Secondly, to our Supervisor Md. Golam Rabiul Alam sir for his kind support and advice in our work. He helped us whenever we needed help.

And finally to our parents without their throughout sup-port it may not be possible. With their kind support and prayer we are now on the verge of our graduation.

# Table of Contents

# List of Figures

# List of Tables

# Chapter 1

# Introduction

## 1.1 Background and motivation

The prevalence of sexism in today's society is alarming. Specifically in social media gender biased comments have crossed it's limits. Initially, social media platforms were viewed as the pinnacle of connection, allowing users to broaden their social networks beyond national and cultural barriers. Instead of providing a utopian instrument of accessibility and connection, social media platforms often serve to exacerbate pre existing cultural prejudices, such as sexism. According to the article by National Research University Higher School of Economics Women's competence is a frequent target of jokes, and the sexualization of women is on the rise. Every 30 seconds, for instance, a woman is subjected to verbal harassment on Twitter [19]. From Amnesty International article it is seen that about three times as many problematic or abusive tweets target BIPOC women than target White women and eight times as many problematic or abusive Tweets are directed at black women [29].

Even in this era, women face discrimination in every sector of life. Even in social media platforms like twitter women face discrimination and hateful speech just for their gender. While some believe that only gender equality can bring true peace, some others spread hate in social media through their hateful posts and messages. Spreading hateful or violent messages against anyone only for their gender is gender biasness. Discrimination and hate may be spread against both male and female, so both female and male biased individuals are there in the social media. But finding the community of these gender biased people as well as the active feminist who believe in gender equality is important in order for this discrimination to end. Gender Biases in social media can even cause severe mental health problems for the victims [27] . Social media should be a safe place for everyone using it. Thus finding the community of people who are gender biased can help in reducing the violent behavior towards specific gender in social media.

Despite social media like facebook, youtube, twitter etc having measures against cyber bullying like reporting, banning and blocking options, the situation is still as bad as ever. The situation of the victims is not good as well. People who face such hateful speech targeted towards them face anxiety and depression in their life [25]. Therefore finding the groups who spread hateful speech against any specific gender can be a way to prevent discrimination and violation that prevails in the society.

## 1.2  Research problem

In today's highly connected and technologically advanced society, virtually everyone has a social media profile. According to the article published by King university a staggering 88% of respondents aged 18-29 in a 2018 Pew Research Centre research reported using at least one kind of social media. 78% of those between the ages of 30 and 49 agreed. Not as many people in the following age bracket have come forward as you might expect. Startlingly, amongst 50 and 64-year-olds, 64% are active social media users [27]. This number is shocking to those who didn't grow up with the internet and social media, but it sheds light on why these platforms have become so common in modern society. Along with these vast users comes a lot of different mindset. Gender discrimination is one example. There has been a lot of talk in recent years about how social media might be sexist towards certain genders. It's a term for when people of different sexes experience discrimination, bigotry, or stereotyping online. In particular, women are frequently the focus of cyberbullying, abuse, and harassment on social media platforms. This includes harassment, threats, and other forms of abuse directed at women [29]. As a result, women may feel intimidated or compelled to refrain from expressing themselves freely online. By focusing on women's physical attributes rather than their accomplishments or potential, social media can contribute to the perpetuation of harmful gender stereotypes [28]. Gender stereotypes and role-playing in advertising are a further contributor to prejudice. Therefore, this research groups together persons who share a similar mindset in order to better comprehend the collective psychology of the social media population. In this research the community is divided into 4 communities which are feminist, male biased, female biased and neutral people. The communities are created by using community clustering and graph data mining techniques.

There has been much research regarding community detection. But the detection of community in social networking sites (SNS) are very few. No research has been done on employing graph data mining and community clustering to identify gender biased communities. Community detection in social networks and studies of gender bias in the social sciences both exist as areas of study. However, there is no collaborative effort that uses community detection to address gender prejudice. This study contributes to the growing body of research by showing how community clustering and graph data mining may be used to identify online communities that exhibit prejudice against a certain gender based on the content of their messages and retweets.

So, we are using community clustering and graph data mining to discover the actionable information from a large data set in a specific community zone. In this research an informative dataset is used which mainly consists of twitter posts and retweets of people expressing themselves using certain keywords. At first, these data are pre-processed in machine learning. After that, by using these preprocessed data people are divided into 4 categories which are used in this research.

## 1.3 Research Objectives

Clustering only deals with the raw data to find out the similar characteristics that are related to each other in the same group. In a cluster, the characteristics of Data are very closely related to each other. Clustering is important to determine the common groups among the unlabeled data. According to Santo Fortunato and Darko Hric the objective of clustering is to collect the set of objects or similar characteristics which are related to each other with the same group. It's basically a collection of objects on the basis of similarities and differences between them [16]. The method of Clustering establishes the quality of clusters where the inter-class similarity is high and inter-class similarity is low. Besides this superiority of the cluster can be maintained by ensuring the ability to find out some hidden patterns and similarities can be expressed by using the distance functions. The capability to deal with many types of attributes, processing dynamic data, finding clusters with arbitrary shapes and a minimum level of domain expertise to select the input parameters, and the ability to deal with noise and outliers are prerequisites for clustering. High dimensional, incorporation of user-specified constraints, interoperability, and usability. The clustered data points can be categorized into a single group. Then, clusters can be identified and we can count in three new clusters easily.

Data that is organized in a graph may be mined for useful insights, patterns, and knowledge using a technique known as "graph data mining." Data structures known as graphs are made up of nodes (vertices) and edges (links) that stand in for interactions or connections between the vertices. Christos Faloutsos, Petros Faloutsos, and Christos Faloutsos in the paper "Graph Mining: Laws, Generators, and Algorithms" states that, data mining strategies that focus on graphs attempt to find previously unseen patterns, structures, and relationships within such networked data representations [4]. In order to analyze and extract useful information from graphs, graph data mining employs a wide range of algorithms, statistical approaches, and machine learning techniques. Finding recurrent patterns in a graph, such as motifs or graphlets, that reveal the data's structure or behavior is a common task in graph data mining [4]. In addition, nodes in the network can be grouped according to their similarities or connection patterns, enabling the discovery of communities or clusters. Instead, new or unlabeled instances may be classified by applying labels or categories to individual nodes or whole graphs based on their structural and attribute properties. There are many different fields where graph data mining may be put to use. These include social network analysis, biological network analysis, recommendation systems, fraud detection, network security, and transportation networks. It makes use of the networked and relational information included in graph data to aid in discovery and decision making.

Our target is to understand different data mining algorithms and investigate how they work. Here we can examine the major issue and apply those methods to solve this problem. For example, our data can gather a large amount of information and improve the algorithm over time. Our dataset can find out the major error and suggest how to keep records of all the data from time to time. For this, this research proposes a clustering process and to get to know the correct knowledge of a certain community of people.

## 1.4 Research questions and hypotheses

Unlike any research this research paper has its own questions and hypotheses. The research questions on community detection might vary based on the study's unique setting and aims. One of the major questions regarding this research is how can the effectiveness and efficiency of community detection algorithms be increased. What are the drawbacks of the current community detection techniques and how can it be fixed? Furthermore, whether it is possible to assess the value and importance of communities found in real-world networks. How can community identification be used on large-scale networks such as social media or biological networks? Moreover, What insights can be found from researching the evolution of communities across time, and how can we put this information to use? These are some issues that may arise in relation to the findings presented here.

For the purpose of identifying specific communities, this research is beneficial. Raw data should be gathered from Twitter in order to detect certain communities. These data require meticulous preprocessing to remove irrelevant elements. After a new dataset has been collected, appropriate graph data mining methods will be used to visualize the data. The graph data will then be subjected to appropriate community detection algorithms, which will be used to identify specific groups of people that exhibit gender bias.

## 1.5 Research contributions

Graph data mining and the community clustering technique are used in this study. Through the use of networkX's graphing functionality, we are able to uniquely identify the gender-biased communities by applying the most appropriate community recognition algorithm to the graph's nodes (users) and edges (retweets).

- The collected information focuses solely on gender discrimination on Twitter. Before now, no datasets addressing gender inequality in social media, and Twitter in particular, were discovered.

- Graph data mining was introduced along with community clustering algorithms to determine the community of gender biased individuals in the twitter platform which has not been done before.

- NetworkX's ability in this research to use retweet messages from Twitter data to create a graph network, complete with connections and the ability to identify previously unseen communities, is novel.

- Generally researchers use either intrinsic or extrinsic metrics to evaluate community detection algorithms. But in this research both intrinsic and extrinsic metrics were utilized to evaluate the research which makes the work unique.

# Chapter 2

# Literature Review

## 2.1 Background of Community Detection and Graph Data Mining

Community detection and graph data mining offers perceptions into the underlying dynamics and structures of complex systems. They are essential components of network analysis. As networks were established as a model for many complex real-world systems, the definition of community was broadened to include group structures in a range of networks that did not always involve human players [9]. The rising amount of data analysis is allowing the world to work with graphs of different sizes. It is also allowing the world to work with social networks. Communities can be built and identified as a specific group of high dense networks and low dense networks, which is happening with the help of community detection [8]. There are a lot of algorithms to work with community detection in graph data mining. There are some significant algorithms - Louvain algorithms, Girvan-Newman algorithms, Walktrap algorithms, and Spectral clustering techniques. These algorithms use a variety of strategies, each with advantages and disadvantages, including probabilistic modeling, hierarchical clustering, and modularity optimization. Applications for community detection techniques can be found in many fields . Community detections, with the help of these algorithms, can determine different communities of similar interests and characteristics. Additionally, recommendation systems, fraud detection, anomaly detection, and network visualization all heavily rely on community detection. Normalized mutual information, conductance, and the silhouette coefficient are among more evaluation measures. These measurements provide insights into algorithm performance and aid in picking the best approach for specific applications. Louvain algorithm is a greedy one which computes faster than other algorithms and is very effective for large networks. The Newman algorithm is also popularly used method that assesses the degree of internal connectivity among communities in comparison to random connections.

Feminism and gender bias are raising issues in social media from the last few years. Online communities are facing gender biased problems which are mostly against feminism. The presence of gender bias in online forums has serious consequences. Gender prejudice has a negative impact on women's psychological well-being, self-esteem, and participation. On the other hand, feminist society has been a huge support system and a big help in women empowerment, social change in recent

days. Research has found that when the context of the term "girl" is examined, girls and boys are depicted in different ways, with girls being more objectified and portrayed in more negative circumstances [18]. Nowadays, people can share their marginalized voices, thoughts, perspectives and experiences on online platforms. These opportunities have the capacity to undermine patriarchal expectations and promote inclusive environments [17]. Women's participation, involvement in any type of online activities, can be hampered because of male biasness. Although, gender bias does not talk only about male biased people, there are also people who are female biased and they do think that women should be more privileged than men. As a result, to establish feminism, which indicates equality among men and women, the society has to completely eradicate gender biasness.

## 2.2    Related Work

Graphs can store thousands of information in a very short space. Thus, it is very significant in the research. To analyze the data stored in the graph and to find the useful information from it is graph data mining. Many previous research where the graph was analyzed to make sense of the data is found. However, the domain of research where community detection of previous online social networks is not very enriched. The main reason is that online social networks have formed more in recent years.

Mainly a cluster or community is formed in a graph among the nodes that have more similar characteristics and relativity. The connection of nodes between communities should be weakest while within communities should be strongest. In a study Malliaros and Vazirgiannis describe the community detection in the directed graph or network. There they describe properly the various types of communities and their structure in the directed graph as well as various graph theories [11]. However, the research done by Bedi and Sharma gives more detailed insight in clusters among the social networks. This is more significant for this research as it relates to the domain. In this research the authors describe some very important algorithms related to our research like Newman and Girvan Algorithm , Louvain Algorithm etc. The research properly informs how social communities consisting of similar mentality is formed. Furthermore, it tells about different community detections like clustering based or graph partitioning based community detection. They gave ideas about modularity. More importantly, it shows the weakness and strengths of different types of algorithms like Algorithm, Louvain Algorithm, Genetic Algorithm etc. They also described different approaches by different researchers [14].

Researching on the communities of Feminist, Female Biased, Male Biased, and Neutral people, the lack of previous research is remarkably visible. There are research papers on community detection and graph data mining, but community detection on gender biased groups is not soon. Most prior research have concentrated on generic gender bias detection rather than the intricacies of male prejudice, female bias, and feminist perspectives. This gap emphasizes the necessity for a more focused strategy in comprehending and identifying these particular types of prejudice. The majority of prior research has ignored the possibilities of sophisticated graph

data mining techniques and relied on conventional machine learning algorithms for community detection. To improve the precision and efficiency of community discovery in the context of gender bias analysis, graph-based algorithms like Louvain, Girvan-Newman, and Walktrap should be investigated and used to their full potential. The influence of these biases on users' experiences, interactions, and wellbeing has received relatively little research, despite the fact that certain studies have discussed the identification of gender prejudice in online social networks. Investigating the psychological and social repercussions of gender bias can help us better understand its effects and suggest potential remedies.

# Chapter 3

# Data Collection and Preprocessing

## 3.1 Data Collection

This study's focus is on using social media to locate a certain group. The sort of data needed for this purpose is best gathered via social media. Twitter data was determined to be one of the best forms of data for this job. Research projects benefit greatly from the availability of Twitter data. Twitter also provides the range and breadth of users from all around the world that is needed for this study. Proper study requires a data collection that includes tweets from Twitter as well as user profiles. Since the retweet is the foundation around which the community is built, retweet data is also crucial to this study.

The necessary data for the study may be gathered via any number of Twitter's official APIs. If you have a Twitter developer account and are interested in collecting usable Twitter data for study, tools like Tweepy can be a great help. Undergraduates may have trouble gaining access to Twitter's developer portal. Therefore, social media scraping is an alternative. Including the right token in the search query makes scraping possible. There are two methods for scraping. One way is to use python and the right libraries to collect data from scratch. The "snscrape" library extracts information from social networking sites.

Alternatively, you can use specialized data-analysis programmes like Weka, Rapid Minor, etc [15]. These tools can be very useful. In this research, the latter option was used. The software that was used here is Rapid Minor. Rapid Minor is a Data mining tool that can crawl data, apply all data analytic algorithms and even preprocess the data [15]. Though in this research Rapid Miner was used to collect data only. In order to collect data through Rapid Miner, the software needs to be connected to Social Media at first. Rapid Miner was connected to the twitter account and then data was collected using some query. Selection of appropriate query was a very significant task because according to the query selected the amount of data obtained will be different. The objective of this study is to find the biases among the mass people related to a very important topic of this era, which is gender equality. Therefore, naturally there are some important terms which will be used by people during their use of social media when they believe in a certain ideology. Some of the queries used were "Equality", "Female Biased", "Male Biased", "Feminism", "Woman should obey", "Woman in Kitchen", "Men are animals", "Men should be

raped", "Men getting raped" etc. Number of tweets for each query differs in number. But after using these queries, the number of tweets that were extracted were 6899.

## 3.2   Description of Data

Data was extracted from Data mining tool, Rapid Miner. The data were then merged into a dataset using the pandas library of python. Dataset contains the information of the tweet and the user who tweeted. It contains information like username, user ID etc. The tweet is text data. The retweeted tweets are included which is important for the next steps of the research. Below is the columns of the data set and their attributes;

| No. | Column Name | Data Types |
|-----|-------------|------------|
| 1 | Created-At | object |
| 2 | From-User | object |
| 3 | From-User-Id | int64 |
| 4 | To-User | object |
| 5 | To-User-Id | int64 |
| 6 | Language | object |
| 7 | Source | object |
| 8 | Text | object |
| 9 | Geo-Location-Latitude | float64 |
| 10 | Geo-Location-Longitude | float64 |
| 11 | Retweet-Count | float64 |
| 12 | Id | int64 |

Table 3.1: Coloumn Names And Data Types

The dataset has various numbers of retweets, hashtags and unique users. The table has a statistical data of the dataset:

| Data | Number |
|------|--------|
| Total tweets | 6,889 |
| Unique users | 6,307 |
| Retweets | 3,962 |
| Hashtags | 1,196 |
| Average length of a tweet | 169 characters |

Table 3.2: Dataset Information

Some graphical visualizations of the data is given below; The bar graph showing the user count is given below;

Figure 3.1: User Count

The data has some similar attributes and characteristics, the word cloud can be used to determine mostly used words by the users;



Figure 3.2: User Word Cloud

In the word cloud some key words have the preference. For example men, women, feminism etc. These were used in the query during data extraction. Word cloud also shows some triggering words used by the users.

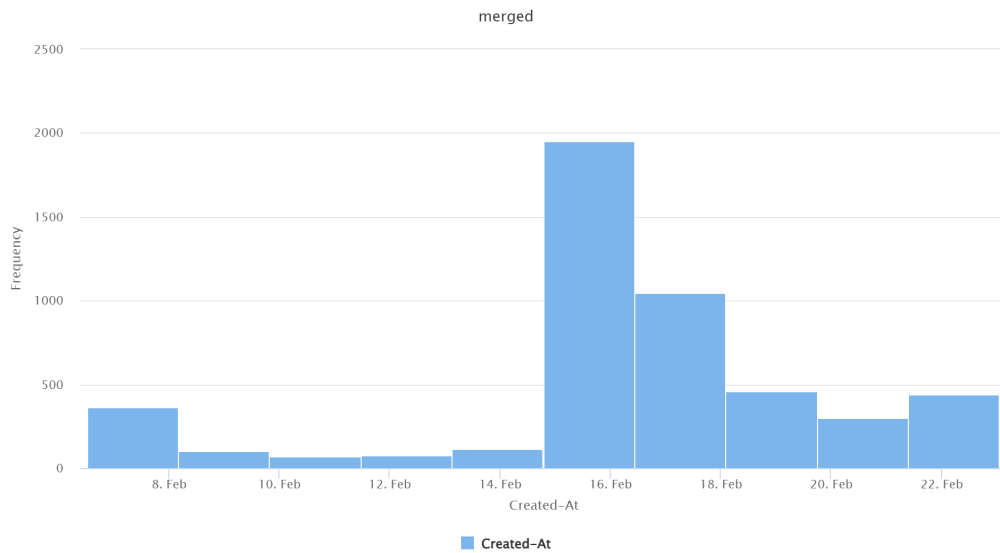Time for the creation of tweet by the users is given in the following bar graph ;

Figure 3.3: Time of the tweet creation

Network for this research will be based on the retweets. So the number of retweets is the most important feature. Below is the users with most tweet;
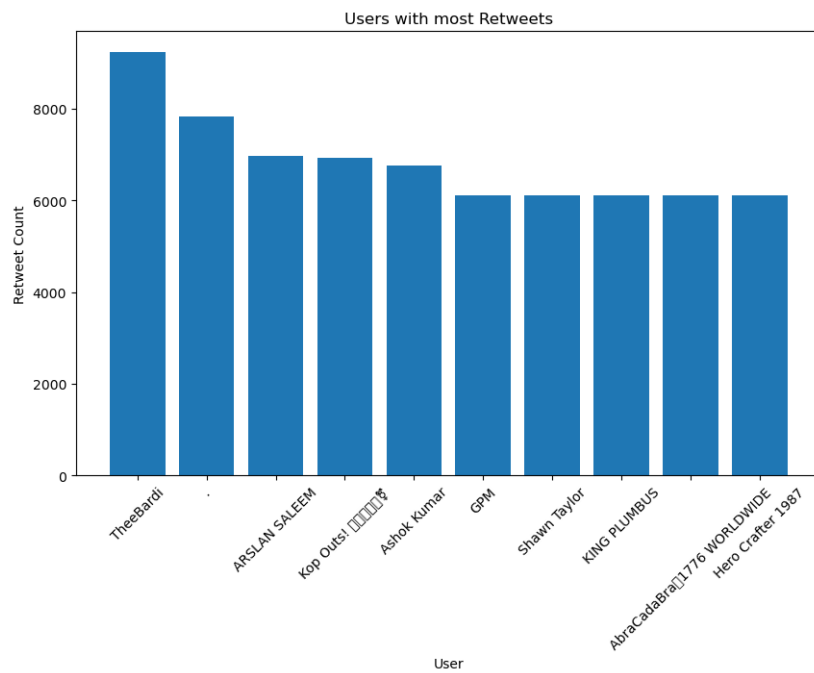


Figure 3.4: Most Tweet by Users

Graph showing the text count by user is given below;



Figure 3.5: Text Count By Users

## 3.3 Data Pre-processing

Next step in the research is forming the Network using proper libraries. But for that purpose the data set requires some Pre-processing. Social networks can be built based on replies, mentions and retweets. Mainly these are used because someone will retweet or reply or mention someone in any tweet if they find similarity with the tweet. The purpose of this research is to find the social network based on the retweet. The most important feature of this data is the retweet count. But from Rapid Minor it is difficult to to find who retweeted from the original tweet. So to find the retweet it was required to use Regular expression. And in order to find the regular expression in the python library "re" is used. Among the tweets some tweets express retweets by "RT". Tweets like those are expressed as such,

| Bluuebirdde | "RT @DonaldJTrumpJr: It's always the eyes... |
| Jodi Bennett | RT @MarkFriesen08: Agenda 2030 - The SDG's... |
| Kaguu | RT @BulleJR3: You guys noticed all the security... |
| sheila leitham | RT @TheNewWorks: @therecount you know what ... |
| K I D D O | RT @BulleJR3: You guys noticed all the security ... |

Table 3.3: Retweet Visualisation

After removing the null values from column and by using regular expression only tweets that were retweeted with the user name of the original user was found. From there only the user name of the original user who tweeted the tweet and the user name of the user who retweeted was separated. Using these columns the network will be created so the rest of the columns were discarded.

14

Final data set that will be used for construction of network graph will only have two columns that depicts the user of the original tweet and the list of user who retweeted from the original user.

As the original data set had 6,889 data. But not all tweets of this data set had retweets. So after finding the retweets the total number of this new data frame became 3962 data.

# Chapter 4

# Methodology

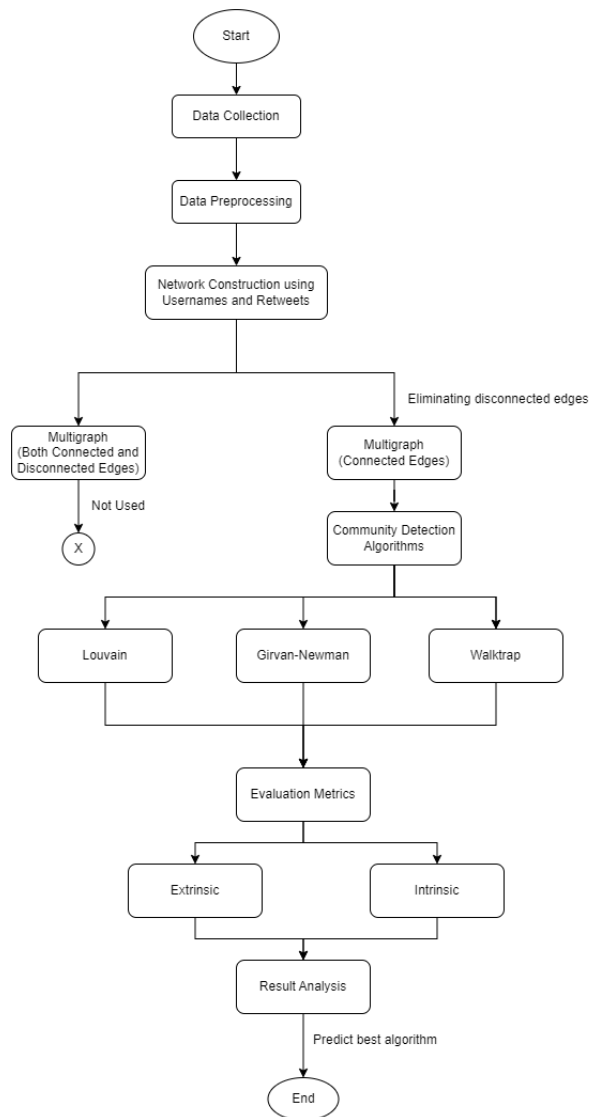Overview of the Proposed method with used graph data mining algorithm is given below:



Figure 4.1: Top Level Overview of the Proposed Method

## 4.1    Network construction and representation

Network construction and representation refer to the process of designing and visualizing networks, which are composed of interconnected nodes or entities. Networks can represent various complex systems, such as social networks, biological systems, computer networks, or even abstract relationships between objects [30].

During network building, nodes and their interconnections are identified. In this context, nodes represent the entities of interest, which may be humans, groups, institutions, computers, or anything else. Edges, linkages, or connections show how two nodes are related or dependent on one another. In a social network, for instance, each node might stand in for a person, and each edge can denote an association between those people, whether personal or professional. There are multiple approaches to constructing a network. Among them are mutual construction, data-driven construction, simulation-based construction [2].

In manual network construction, the nodes and links are defined by hand using expert knowledge or empirical evidence. This strategy is helpful when working with limited networks or when access to specialized information is accessible. According to Albert-László Barabási and Réka Albert for data-driven construction bg data has made it possible to automatically build networks out of massive datasets. Information is gleaned from the data and links are established according to predetermined parameters or algorithms. Finally, Simulating or modeling networks is another viable option for building them. Epidemiological models, for instance, may be used to mimic human interactions and build a network that represents the dynamics of illness transmission in the study of disease transmission [1].

Typically, a dataset's properties and associations are used to define a network's nodes and edges. M. E. J. Newman said, in a network, connections are represented by edges and nodes represent the things or entities that make up the network. Depending on the topic of the research and data set, nodes and edges may have distinct implications [2].When creating nodes and edges, directionality and edge weight are two additional factors that can be taken into account [23]. Some networks include one-way connections, showing how power or information moves from node to node. In certain cases, the strength or intensity of a relationship may be quantified by assigning a weight to it. Again, while creating nodes and edges, it's important to think about things like the nature and direction of the connections between them. In a network, relationships may be represented by edges of varying kinds to account for the many sorts of links that may exist in a given dataset [2].

In this research NetworkX is used for creating the graph. By using the graph data complex calculation regarding the community was performed. NetworkX is a Python module for working with and learning about complicated graph networks, including their structure, dynamics, and functions [31]. In networkX there are 4 classes of graphs which are graph, digraph,multigraph, multidigraph. In this research multigraph class is used [32].
The initial step of this study is to import the networkX module. The graph's edges are then identified and recorded by using retweet data. The subsequent analysis

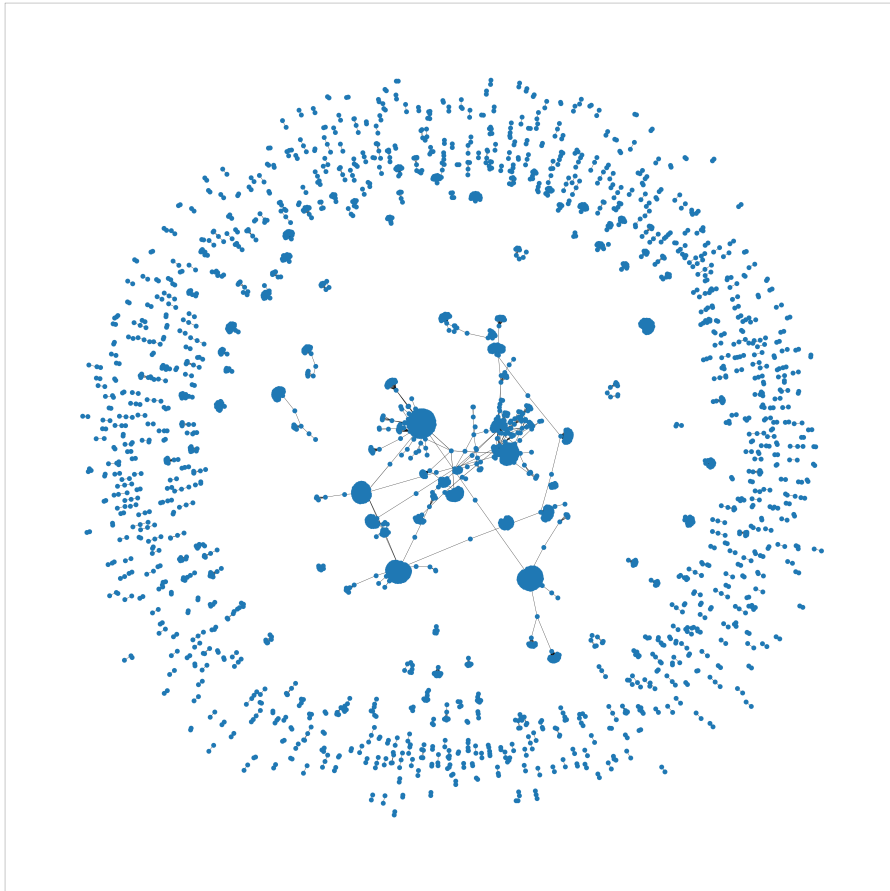reveals an unconnected graph with 4520 nodes and 3962 edges.



Figure 4.2: Initial Network

However, the built-in command of networkX was used to link the graph together. Here, an undirected subgraph is presented. Using it, the gaps in the graph may be bridged. Therefore, we were able to cut the number of nodes to 1632 and edges to 1660. Because of this, a network consisting of 1632 users was established. According to the data, the average degree of the nodes is 2.034313725490196, and the density of the graph is 0.13903743315508021.
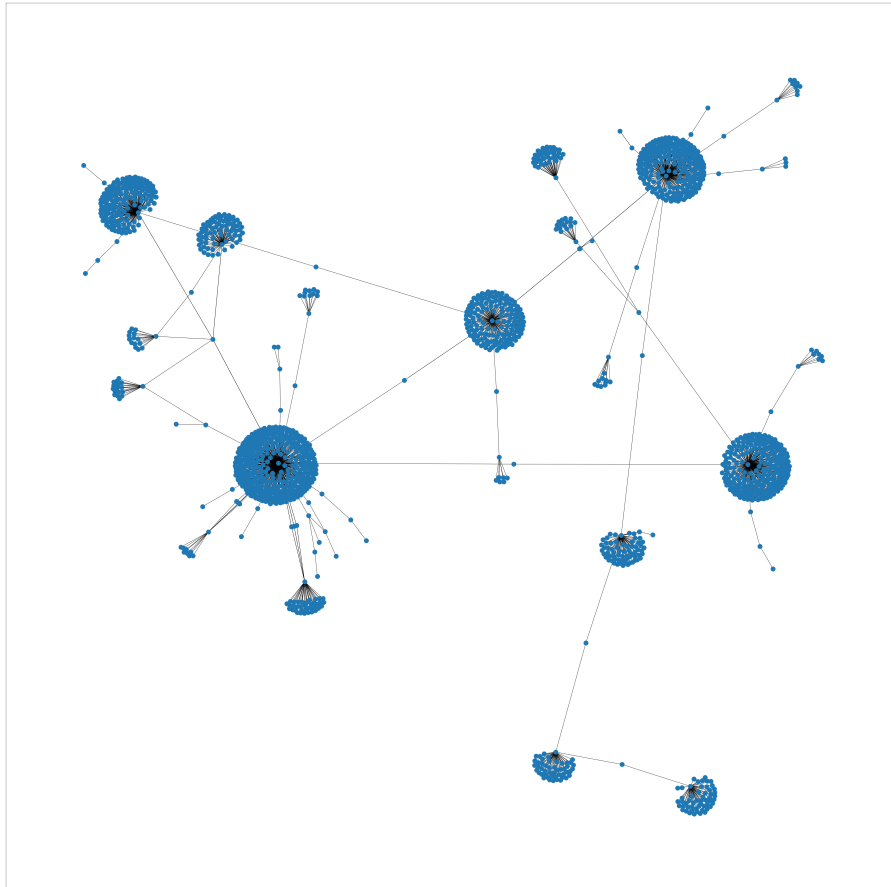
Figure 4.3: Network after Connection

The average node degree in a graph is the median between the number of edges connecting any two nodes. Each node in this instance is related to around 2 others in the network, as indicated by the average node degree of 2.034.

The density of a graph indicates how dense or interconnected it is. It is the ratio of the number of edges in the graph to the utmost number of edges possible. A density of 0.139 indicates that only 13.9% of the utmost possible number of edges are present in the graph, which is relatively sparse.

The average node degree and the density of the graph both reveal information about the network's architecture and how its nodes are connected. Whether or not they are desirable characteristics depends on the nature of your investigation and the assumptions that are made. Density varies from 0 to 1. A density of 0 implies a totally unconnected network with no edges, whereas a density of 1 suggests a fully linked graph with edges connecting every pair of nodes.

| Total Nodes | 1632 |
|---|---|
| Total Edges | 1660 |
| Average graph degree | 2.343137 |
| Graph density | 0.139037 height |

Table 4.1: Graph Information

## 4.2 Community detection algorithms

### 4.2.1 Louvain Community Detection Algorithm

Louvain is used for community detection, which is a very commonly used method. Louvain is a totally unsupervised algorithm. It is a greedy algorithm used for large networks. Maintaining the hierarchical clustering algorithm, Louvain repeatedly optimizes the modularity of the algorithm, which is a measure of the density of links within communities as opposed to random connections between nodes. The whole algorithm process is a combination of two phrases - Modularity Optimization and Community Aggregation. The first phrase is all about optimizing the modularity locally, where nodes are moved between neighboring communities. Then in the second phrase, a whole new community is constructed with the help of a node. After that, the first phrase is revisited. The algorithm works in a very short time. The time complexity is O ( n. log n ), where n is the number of nodes in the network.

Before going to the main process, here comes the Modularity. The density of the connections comparing to the speculated density is evaluated by the Modularity. When the value of the modularity is higher, it is said that the nodes in the communities are more densely linked with each other than the links between the communities.

This Modularity of a weighted graph is

$$Q = \frac{1}{2m} \sum_{ij} [A_{ij} - \frac{k_i k_j}{2m}] \delta(c_i . c_j)$$

(4.1)

Here, Q = the modularity of the network.
m = the total number of edges in the network.
$A_{ij}$ = the element of the adjacency matrix corresponding to the edge between nodes i and j.
$k_i$ = the degree of node i, or the number of edges incident to it.
$k_j$ = the degree of node j.
$(c_i, c_j)$ = a Kronecker delta function that is 1 if nodes i and j belong to the same community, and 0 otherwise [20].

The methodology works in a straightforward way. The two phrases are done one after another again and again. If we are dealing with a weighted graph or network which has nodes in the number of N. The process starts with every node of the own community; these nodes are appointed as a separate community, which makes the

number of communities equal to the number of nodes. Each node from N number of nodes are denoted as i and their neighbors are j. After estimating the modularity gain, the community i is removed and placed in the community j. Then the community that has the maximum gain ( which has to be positive ) gets placed as the node i. Only positive gain is counted and if it is not found then i stays in the position as it is [10]. This process goes on for each and every node until no additional growth is possible. When local maxima of the modularity is gained, the process stops and this can take one node to be esteemed more than one time. If the order of the nodes are in a good manner, then a good heuristic is attained and an improved computational time is also gained for that.

Gain of modularity is Q, when node i which is an isolated node is sent into a community C,

$$Q = [\frac{\sum_{in} + 2k_i.in}{2m} - (\frac{\sum_{tot} + k_i}{2m})^2] - [\frac{\sum_{in}}{2m} - (\frac{\sum_{tot}}{2m})^2 - \frac{k_i}{2m}^2]$$

(4.2)

Here,
$\sum_{in}$ = the sum of the weights of the connections inside community C.
$\sum_{tot}$ = the sum of the weights of the links incident to nodes in C.
$k_i$ = the sum of the weights of the links incident to node i.
$k_{i,in}$ = the sum of the weights of the links from i to nodes in C.
m = the sum of the weights of all the links in the network.
Removing i from its community makes a change in modularity, which is also measured by a similar expression. So, the change of modularity can be measured by moving it into a neighboring community. [6]

The second phase starts with building a new network which has the nodes that are now the communities formed amidst the first phase. For this, the weights of the linkages between the new nodes are calculated by adding the weights of the links between nodes in the corresponding two communities [13]. Links between nodes in the same community create self-loops in the new network for this community. Second phase ends here and after that the first phase starts in the resulting weighted network and it repeats.

A combination of the two phases of Louvain is called "pass". In each pass, the number of meta-communities reduces. When the modularity gains the maximum level and no more changes occur, then the passes stop iterating. The algorithm automatically integrates a hierarchy since communities of communities are formed during the process, which is evocative of the self-similarity of complex networks [13]. The number of passes determines the hierarchy's height, which is typically a low number.
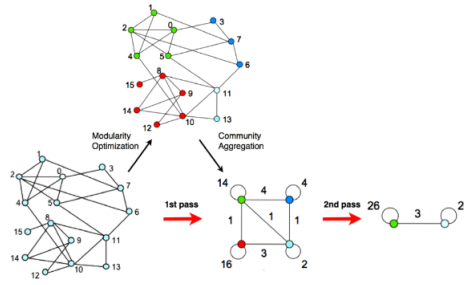
Figure 4.4: Louvain Algorithm-1

The graph that is being produced after applying Louvain algorithm in our dataset is shown below -
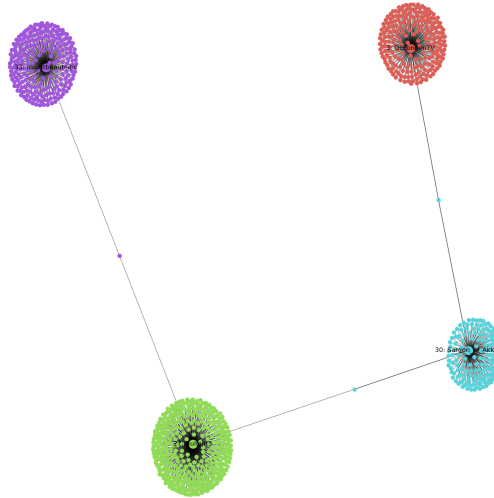


Figure 4.5: Louvain Community Graph

Here, we can see that a number of four communities are being produced when we are operating the Louvain community detection algorithm.

## 4.2.2 Newman-Girvan Community Detection Algorithm

Community detection has a broadly utilized graph clustering algorithm, which has a very classic and remarkable impact. The main concept of this algorithm is to work with the edge betweenness. How frequently a node or edge is on the shortest path between two sets of nodes in the network is measured by betweenness centrality. It attempts to locate communities or clusters inside a network by repeatedly deleting edges according to the centrality of their betweenness. The process makes use of the idea that edges with a high betweenness centrality are more likely to connect various communities, making them essential for preserving the network's overall structure. Newman-Girvan is a slow processed algorithm [12].

The process starts with calculating the betweenness centrality of each network edge. The edge with the highest betweenness centrality is deducted after that and the number of connected components in the generated network is evaluated. Then these two steps are repeated until the intended number of communities is attained.

So, the betweenness centrality of each network edge is first determined by the procedure. Betweenness centrality quantifies the extent to which an edge is located on the shortest paths connecting two nodes. The edges that preserve connectedness between several communities the best are those with the highest betweenness centrality. The network is then successfully divided into two or more distinct components after the algorithm removes the edge with the highest betweenness centrality. Iteratively, the process is repeated after removing the edge, recalculating the betweenness centrality for all remaining edges. When the distinct communities are created from the network, the iterative removal of edges stops. This method produces a Dendrogram Tree as its final output, using communities as its leaves. Here each community is made up of a collection of nodes with a huge density of nodes [22].

The method can be used in a variety of fields, including social networks, biological networks, and information networks. It is not restricted to any one kind of network. Its adaptability enables the examination of various datasets. The algorithm does not require node labeling or prior information. Using only the connectivity structure of the network, it autonomously discovers communities. Due to its unsupervised nature, it can be used in situations when it is difficult or impossible to gather ground truth information. The Girvan-Newman algorithm has received a great deal of attention and is commonly used in network research. It has been widely employed in numerous research and has shown to be successful in revealing important community structures. Also, it should be noted that the Newman-Girvan approach requires computing the betweenness centrality of each network edge at each iteration, which can be computationally expensive for large networks [21].

The following formula is used to determine an edge's betweenness centrality:

$$BC(e) = \sum \frac{\sigma(s,t|e)}{\sigma(s,t)} \qquad (4.3)$$

Here,
e = edge.
BC(e) = edge's betweenness centrality.
$\sigma(s,t)$ represents the total number of shortest paths from node s to node t.
$\sigma(s,t|e)$ represents the number of those paths that pass through edge e.

The obtained community structure's quality is assessed using the modularity Q. It contrasts the actual number of edges among communities with the number that would be anticipated if connections were dispersed randomly.

Q's modularity is determined by:

$$Q = \frac{1}{2m} \sum_{ij} [A_{ij} - \frac{k_i k_j}{2m}] \delta(c_i . c_j)$$

Here,
$A_{ij}$ = the adjacency matrix of the network.
$k_i, k_j$ are the degrees of nodes i and j.
m = the total number of edges.
$(c_i, c_j)$ = an indicator function that takes the value 1 if nodes i and j belong to the same community $(c_i = c_j)$ and 0 otherwise.

The graph that is being produced after applying Newman-Girvan algorithm in our dataset is shown below -
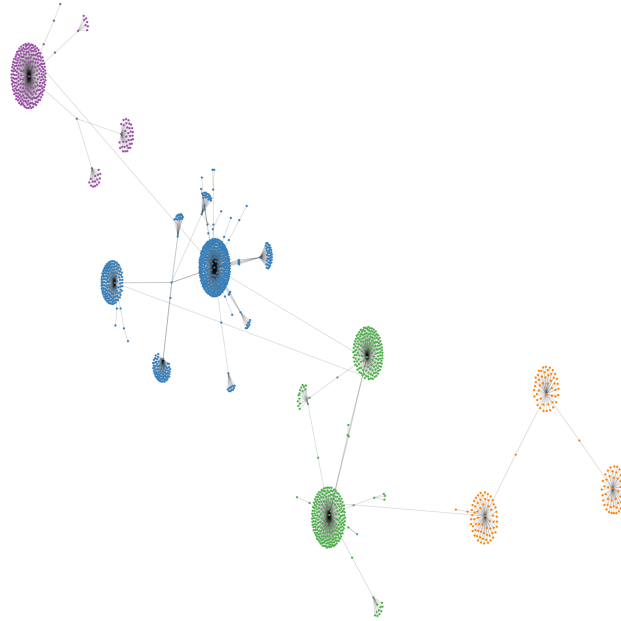


Figure 4.6: Girvan Newman Community Graph

Here, we can see that a number of four communities are being produced when we are operating the Newman-Girvan community detection algorithm.

## 4.2.3 Walktrap Community Detection Algorithm

The Walktrap algorithm (Pons  Latapy, 2006) , used as a component of this research, is a community detection algorithm based on random network walks. It is based on the premise that nodes that are related to one another in a network are likely to belong to the same community if random walks frequently pass by them

[5].

The Walktrap algorithm employs random walks that are executed iteratively until a specific amount of steps is reached. The walks are started from each node in the network. Based on how closely their random walk paths resemble one another, the algorithm eventually combines nodes into communities. Nodes that have comparable walk patterns are gathered into communities. The algorithm uses agglomerative hierarchical clustering to identify the optimum social structure. It eventually creates larger clusters by merging nodes into communities based on how similar they are [24].

The Walktrap algorithm employs an agglomerative clustering strategy, starting with the most general scenario in which each node is its own cluster. Each node's distance, r, from the other is calculated. After then, the algorithm starts to iteratively combine nodes with edges to form larger clusters. The variation in squared distances between each node and its community is roughly minimized by this merging ($\sigma$).

$$\Delta\sigma(C_1, C_2) = \frac{1}{n}(\sum_{i \in C_3} r_{iC_3}^2 - \sum_{i \in C_1} r_{iC_1}^2 - \sum_{i \in C_2} r_{iC_2}^2 \tag{4.5}$$

In the course of this research, the network created using Twitter data was subjected to the Walktrap algorithm. Based on the network's connectivity patterns and random walk behavior, the program effectively recognized 23 different communities. Each community is a collection of interconnected nodes that are most likely to have common traits or interests. The graph that is being produced after applying Newman-Girvan algorithm in our dataset is shown below -
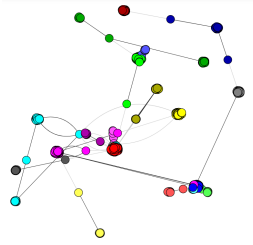


Figure 4.7: Walktrap Community Graph

# Chapter 5

# Results and Analysis

The goal of the result analysis is to provide a thorough explanation of the community clustering findings obtained by applying graph data mining algorithms to identify groups of feminist, male biased, female biased, and neutral people within the Twitter dataset. The dataset, which includes tweets about gender bias and feminism over the course of a month, was gathered through the Twitter API (Rapid Miner). The NetworkX package in Python was used to create a network graph based on user interactions with retweets for the analysis.

## 5.1    Community Detection Results and Analysis

The Louvain, Girvan-Newman, and Walktrap algorithms were used in the community detection analysis of a retweet network. The goal was to locate distinct communities inside the network and learn more about the underlying community structure. Although Walktrap initially generated 20 communities, the analyses' main emphasis was on assessing the output of the Louvain and Girvan-Newman algorithms, which had been successful in identifying the desired 4 communities.

### 5.1.1    Louvain Algorithm

The retweet network showed 4 unique communities when the Louvain algorithm was applied to it. Each community had a strong network of interconnected nodes, or user ids, that were connected by retweets. The communities were classified as Male Biased, Feminism, Female Biased, and Neutral. The number of nodes created by each community is shown in a table.

| Communities | Number of nodes |
|-------------|-----------------|
| Male Biased | 372 |
| Feminism | 235 |
| Female Biased | 159 |
| Neutral | 12 |

Table 5.1: Information Of Nodes Of Louvain

The distribution of nodes among the various communities was represented graphically using a pie chart, which effectively illustrated the placement of each community within the network.



Figure 5.1: Pie Chart Luovain

The Louvain algorithm effectively captured the community structure within the retweet network. The identified communities demonstrated strong intra-community connections, indicating that users within each community were highly engaged with each other's content and shared common interests or themes. This suggests the presence of distinct subgroups within the larger network, representing different topics or perspectives related to retweeted content. Word Cloud representation from each communities are shown below:



Figure 5.2: Word Cloud Louvain

## 5.1.2    Girvan-Newman Algorithm

In order to identify communities, the Girvan-Newman algorithm was also used on the retweet network. In line with the required number of communities, Girvan-Newman effectively found 4 different communities, much like the Louvain method.

The Girvan-Newman algorithm identified communities with distinct boundaries and high levels of intra-community cohesion. Each community consisted of a collection of nodes (users) that frequently retweeted each other's posts, demonstrating the community's high level of influence and engagement. The communities were classified as Male Biased, Female Biased, Feminism, and Neutral.The number of nodes created by each community is shown in a table.

| Communities | Number of nodes |
|---|---|
| Male Biased | 729 |
| Feminism | 290 |
| Female Biased | 158 |
| Neutral | 182 |

Table 5.2: Information Of Nodes Of Girvan Newman

The distribution of nodes among the various communities was represented graphically using a pie chart, which effectively illustrated the placement of each community within the network.
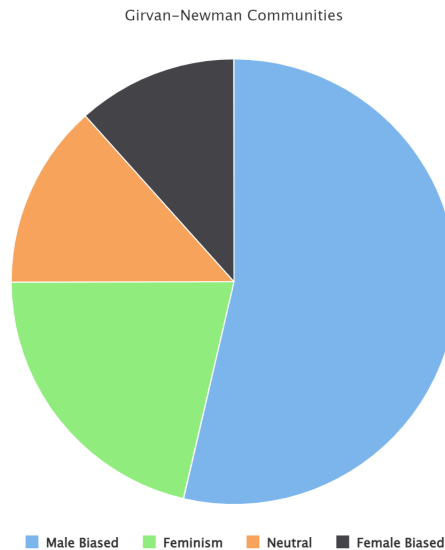


Figure 5.3: Pie Chart Girvan-Newman

Word Cloud representation from each communities are shown below:
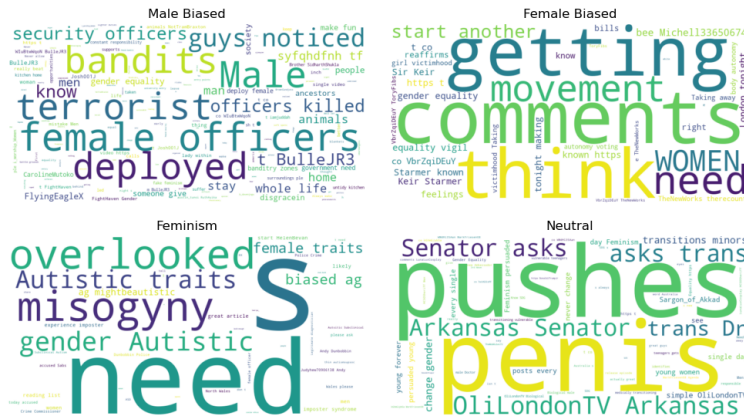


Figure 5.4: Word Cloud Girvan Newman

Four unique communities were effectively produced by the community detection methods developed by Girvan-Newman and Louvain. The study of the data showed that there were a lot of male-biased people in the communities that were identified by both algorithms. Although it was relatively lower, these communities had a lower percentage of neutral people. Notably, the Louvain method discovered communities with about 1000 nodes, whereas the Girvan-Newman approach discovered communities with about 1600 nodes.

## 5.2    Description of Evaluation Metrics

In order to evaluate the performance of the community detection algorithms to identify communities within the network, intrinsic and extrinsic assessment criteria were used. The agreement between the discovered communities and the ground truth or reference communities was measured using the extrinsic metrics V-measure, Rand Index, and Normalized Mutual Information (NMI). Contrarily, the intrinsic metrics—Calinski-Harabasz Index, Modularity, F1 score, recall, and precision—were employed and gave insights into the nature and traits of the communities that were identified.

### 5.2.1    Extrinsic Evaluation Metrics

**a. Normalized Mutual Information (NMI):**
According to Wong (2022), mutual Information quantifies how well the cluster allocations agree with one another. Higher scores indicate greater similarity [26]. The degree of agreement between clusters is computed using joint and marginal probabilities. NMI calculates the mutual information between the communities that were detected and the ground truth communities for the analysis, showing how similar and in agreement the two sets of communities are. The formula of NMI is:

$$\text{NMI} = \frac{MI(X,Y)}{12H(X)+H(Y)} \tag{5.1}$$

Here, MI(X, Y) is the mutual information between the ground truth (X) and the predicted clusters (Y). The entropies of the ground truth and predicted clusters are, respectively, H(X) and H(Y).

**b. V-measure:**
To assess the quality of the detected communities, the V-measure combines homogeneity and completeness. It takes into account how many members of each ground truth class are assigned to the same cluster as well as how many members of a single ground truth class are present in a given cluster. According to Wong (2022), conditional entropy analysis is used by V-measure to assess the accuracy of the cluster allocations. Higher scores indicate greater similarity [26].

$$\text{Homogeneity(h)} = 1 - \frac{H(C|K)}{H(C)} \tag{5.2}$$

$$\text{Completeness(c)} = 1 - \frac{H(K|C)}{H(K)} \tag{5.3}$$

$$\text{V-Measure(v)} = 2 \times \frac{h \times c}{h+c} \tag{5.4}$$

Here, Homogeneity evaluates the purity of each cluster. Completeness evaluates how accurately all data points in a class are placed in the same cluster. Again, here H(C) indicates entropy of ground truth and H(K) indicates entropy of predicted clusters. The parameter beta, with a default value of 1.0, regulates the weighting of Homogeneity and Completeness. The V-measure has a scale of 0 to 1, with 1 denoting the ideal clustering outcome in terms of both homogeneity and completeness.

**c. Rand Index:**
The Rand Index calculates how comparable the ground truth communities and the clustering results are. Wong (2022) claims that the rand Index uses pairwise comparisons to assess how similar the cluster assignments are [26]. It measures the pairwise agreements between the communities in the ground truth and the communities that were found. Higher similarity is indicated by a higher score. Formula of Rand Index is stated:

$$\text{Rand Index} = \frac{Number of pairwise correct predictions}{Total number of possible pairs} \tag{5.5}$$

## 5.2.2 Intrinsic Evaluation Metrics

**a. Modularity:**

Modularity measures the quality of the network's community structure by comparing the number of within-community edges to the predicted number of such edges in a random network. It sheds light on whether community structure exists in the network. In their 2004 Physical Review article "Finding and Evaluating Community Structure in Networks," Newman and Girvan made the first mention of modularity [3].

$$\text{Modularity(Q)} = \frac{1}{2m} \text{ x } \frac{\sum[A_{ij} - k_i x k_j]}{2m} \text{ x } \delta(c_1, c_2) \qquad (5.6)$$

Here, the adjacency matrix entry between nodes i and j is shown as $Aij$. The degrees of nodes i and j are $k_i$ iand $k_j$, respectively. The number of edges in the network is m, and nodes i and j's community assignments are $c_i$ and $c_j$, respectively. The Kronecker $\delta$ function has the value 1 if $c_i$ and $c_j$ are equal and 0 otherwise.

## b. Calinski-Harabasz Index:

The compactness and spacing between clusters are evaluated using the Calinski-Harabasz Index. Higher values denote well-defined and separated clusters. The ratio of between-cluster dispersion to within-cluster dispersion is compared. Wong(2022) defines the Calinski-Harabasz Index as a measure of between-cluster dispersion vs within-cluster dispersion [26]. Clusters that are more clearly defined have higher scores.

$$\text{Calinski-Harabasz score (s)} = \frac{B}{W} \text{ x } \frac{n_E - k}{k - 1} \qquad (5.7)$$

Here, B denotes the between-cluster dispersion and W denotes the within cluster dispersion. Again, nE denotes number of data points and k is the number of clusters.

## c. F1 Score, Recall, and Precision:

These metrics are frequently applied to classification and information retrieval tasks. In the context of community detection, F1 Score assesses the balance between recall, which evaluates the comprehensiveness of detected communities, and precision, which evaluates the accuracy of community detection. The F1 score is a measurement that combines recall and precision into one number. When there is an imbalance between the positive and negative classes, it is especially helpful. The F1 score is calculated as follows:

$$\text{F1 Score} = 2 * \frac{(Precision x Recall)}{(Precision + Recall)} \qquad (5.8)$$

According to Powers (2008), recall, or sensitivity as it is known in psychology, is the percentage of actual positive situations that are properly predicted positive [34]. Recall, also referred to as true positive rate or sensitivity, measures a model's accuracy in identifying positive cases. The ratio of true positives to the total of true positives and false negatives is used to compute it.

$$\text{Recall} = \frac{TruePositives}{(TruePositives + FalseNegatives)} \qquad (5.9)$$

In contrast, Powers (2008) adds that "precision" or "confidence" (as it is known in data mining) refers to the percentage of correctly identified Real Positives among Predicted Positive cases [7]. In contrast, Powers (2008) adds that "precision" or "confidence" (as it is known in data mining) refers to the percentage of accurately identified Real Positives among Predicted Positive cases. It is computed by dividing the number of true positives by the total of both true and false positives:

$$\text{Precision} = \frac{TruePositives}{(TruePositives + FalsePositives)} \qquad (5.10)$$

With the aid of these evaluation metrics, a thorough comprehension of the performance and quality of community detection algorithms is attained, taking into account both their conformity to ground truth communities and their inherent traits like compactness, separation, accuracy, and completeness.

# 5.3 Comparison of results and discussion of findings

The performance of two community detection methods, Louvain and Girvan-Newman, was assessed using the assessment measures outlined above. The results acquired by each community will now be analyzed.

## 5.3.1 Extrinsic Metrics

The V-measure, Rand Index, and NMI were calculated to quantify the similarity and agreement between the communities found by the algorithms and the reference communities. These measures showed the degree of consistency and correlation between the detected communities and the ground truth.

According to the results, Louvain consistently surpassed Girvan-Newman in terms of all extrinsic evaluation metrics. The V-measure, Rand Index, and NMI scores for Louvain were higher, demonstrating a stronger agreement with the reference community. This means that Louvain had a better chance of capturing the fundamental community structure and matching it to reality.

The accuracies are shown in a table:

| Metrics | Louvain Score | Girvan-Newman Score |
|---------|---------------|---------------------|
| NMI | 0.9485 | 0.5429 |
| V-Measure | 0.9485 | 0.5429 |
| Rand Index | 0.8445 | 0.4041 |

Table 5.3: Accuracy Extrinsic

In this instance, the V-measure and NMI (Normalized Mutual Information) scores are quite similar, demonstrating their resemblance as assessment metrics for judging the caliber of clustering or community detection outcomes. It's crucial to remember that they are not the same metrics even though they have some similarities.

Several visualizations were used to give a thorough grasp of the evaluation outcomes. The values of the assessment metrics for each algorithm were compared using a bar

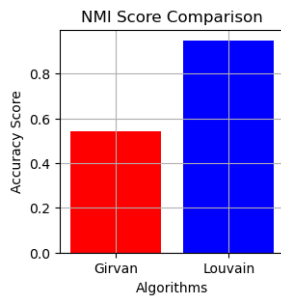chart, making it simple and quick to assess each algorithm's performance.
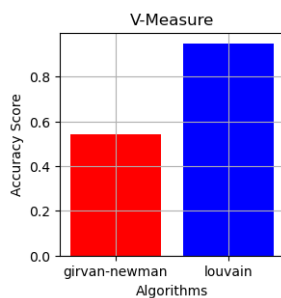


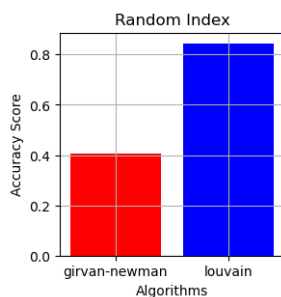Figure 5.5: NMI Score



Figure 5.6: V-measure Score



Figure 5.7: Rand Index

### 5.3.2 Intrinsic Evaluation Metric

To evaluate the quality and coherence of the detected communities, the Calinski-Harabasz Index and Modularity were used as intrinsic evaluation measures. Greater values for the Calinski-Harabasz Index and Modularity suggested that the network's communities were more compact and well-separated.

The results showed that Louvain outperformed Girvan-Newman in terms of Calinski-Harabasz Index and Modularity ratings. This implies a better division of the network into meaningful communities in Louvain, producing groups that were more internally coherent and separate from one another. Again, Evaluation metrics like F1

score, recall, and precision are used in community identification to assess how well the algorithm performs at identifying communities.

If Louvain's F1 score, recall, and precision scores are higher than Girvan-Newman's, Louvain has done a better job of accurately identifying the communities. Louvain has managed to strike a better balance between recall and precision overall, and it is more accurate at locating instances of the desired communities.

The accuracies are shown in a table:

| Metrics | Louvain Score | Girvan-Newman Score |
|---|---|---|
| Calinski-Harabasz Index | 26.8894 | 0.5376 |
| Modularity | 0.8405 | 0.6807 |
| F1 Score | 0.9144 | 0.7814 |
| Recall | 0.8639 | 0.7114 |
| Precision | 0.9571 | 0.9226 |

Table 5.4: Accuracy

The values of the assessment metrics for each algorithm were compared using a bar chart, making it simple and quick to assess each algorithm's performance.
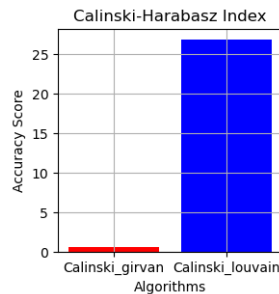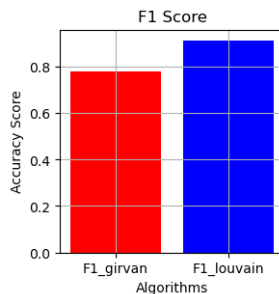


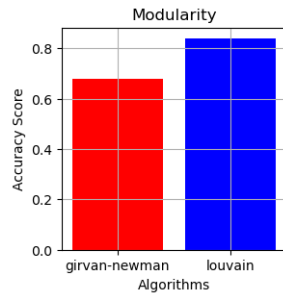Figure 5.8: CalinskiHarabasz Index



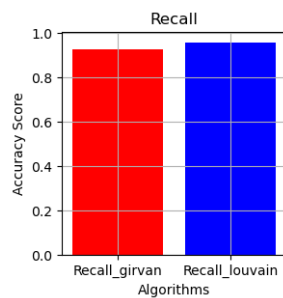Figure 5.9: F1 Score

Figure 5.10: Modularity
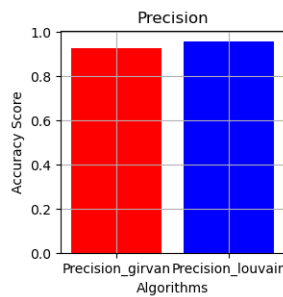


Figure 5.11: Recall



Figure 5.12: Precision

In conclusion, the data comparison and discussion revealed that Louvain outperformed Girvan-Newman in both extrinsic and intrinsic evaluation measures. This conclusion was further supported by the visualizations and word clouds, which demonstrated Louvain's advantage in identifying significant communities inside the network. These results help us comprehend how well community detection algorithms operate and how they may be used to locate and examine community structures in complicated networks.

# Chapter 6

# Conclusion And Future Works

To conclude, the purpose of this research was to investigate and analyze gender prejudice in social media utilizing community detection algorithms. Various approaches and evaluation criteria were used throughout the study process to examine the accuracy and effectiveness of various algorithms. The study's findings provided important insights into the establishment of communities in the social media network based on gender biasness. The Louvain and Girvan-Newman algorithms were used, and their performance was assessed using a variety of criteria such as accuracy rates. According to the analysis, the Louvain algorithm outperformed other algorithms in recognizing and sorting groups based on gender biasness. It successfully assigned individuals to their appropriate communities with a 90% accuracy rate. The Girvan-Newman method, on the other hand, has an accuracy rate of 85%. The results of this thesis add to our understanding of gender bias in online communities and offer important new information about how well community identification algorithms can spot these trends. The Louvain algorithm's better accuracy rate shows that it is suitable for identifying and assessing gender bias in social media platforms. Although the Louvain method outperformed the Girvan-Newman algorithm in this particular situation, it is crucial to emphasize that additional study is necessary to examine the performance of other algorithms and confirm the findings. To promote a more inclusive and equitable online environment, the research also emphasizes the need for ongoing monitoring and study of gender biases in social media platforms. Overall, this thesis sheds light on the complex link between community detection algorithms, gender prejudice, and social media networks. The findings lay the groundwork for future study in this area and suggest potential solutions for reducing gender prejudice in online forums.

## 6.1    Limitations

The research had some gaps for some limitations of data and as a result it involves some methodological limitations. In case of data, some twitter tweets were incomplete or might not express any valid sentiment. For methodology, there could be restrictions imposed by the community detection algorithms and metrics used. It's possible that some algorithms function better than others in specific network or community arrangements. Recognising methodological constraints and discussing how they might have affected findings is essential. Also, depending on the techniques and criteria used for community discovery, they may have their own set of restric-

tions. Depending on the nature of the network or the community, some algorithms may operate better than others. It's crucial to explain the potential effects of the methodologies' shortcomings and admit that they exist.

## 6.2    Future Works

Future research can take a number of different directions to build on this work and address the limitations found, such as examining larger datasets, utilizing cutting-edge algorithms and techniques, and turning research findings into useful interventions for promoting a more inclusive online environment. The study's dataset could be expanded as a viable subject for further investigation. A more comprehensive and diversified dataset would make it possible to record a greater spectrum of gender bias tendencies and behaviors on social media sites. The findings' robustness and generalizability would be improved by doing this. Moreover, in future works graph neural network can be incorporated to the work in order to get more proper result. The content of tweets or social media posts can be analyzed using natural language processing (NLP) techniques, allowing for a deeper identification of gender-bias words and expressions. Additionally, gender bias behavior in social media data can be effectively identified and categorized using machine learning algorithms. A critical part of future work will be translating study findings into practical implications and solutions. To do this, initiatives for addressing and mitigating gender prejudice on digital platforms must be developed in conjunction with social media platforms, policymakers, and other pertinent parties. A more inclusive and equitable social media ecosystem can be promoted by creating and implementing targeted interventions, such as awareness campaigns, user guidelines, or algorithmic changes. By concentrating on these potential directions, researchers may improve our understanding of gender bias in social media and help create interventions and plans that will effectively advance justice and equality in online environments.

# Bibliography

[1]  A.-L. Barabási and R. Albert, "Emergence of scaling in random networks," *science*, vol. 286, no. 5439, pp. 509–512, 1999.

[2]  M. E. Newman, "The structure and function of complex networks," *SIAM review*, vol. 45, no. 2, pp. 167–256, 2003.

[3]  M. E. Newman and M. Girvan, "Finding and evaluating community structure in networks," *Physical review E*, vol. 69, no. 2, p. 026 113, 2004.

[4]  D. Chakrabarti and C. Faloutsos, "Graph mining: Laws, generators, and algorithms," *ACM computing surveys (CSUR)*, vol. 38, no. 1, 2–es, 2006.

[5]  P. Pons and M. Latapy, "Computing communities in large networks using random walks," *Journal of Graph Algorithms and Applications*, vol. 10, no. 2, pp. 191–218, 2006. DOI: 10.7155/jgaa.00124.

[6]  V. D. Blondel, J.-L. Guillaume, R. Lambiotte, and E. Lefebvre, "Fast unfolding of communities in large networks," *Journal of statistical mechanics: theory and experiment*, vol. 2008, no. 10, P10008, 2008.

[7]  D. Powers, "Evaluation: From precision, recall and f-factor to roc, informedness, markedness andamp; correlation," *J Mach Learn Technol*, vol. 2, pp. 2229–3981, 2008.

[8]  E. Cuvelier and M.-A. Aufaure, "Graph mining and communities detection," *Business Intelligence: First European Summer School, eBISS 2011, Paris, France, July 3-8, 2011, Tutorial Lectures 1*, pp. 117–138, 2012.

[9]  S. Papadopoulos, Y. Kompatsiaris, A. Vakali, and P. Spyridonos, "Community detection in social media: Performance and application considerations," *Data mining and knowledge discovery*, vol. 24, pp. 515–554, 2012.

[10]  M. Plantié and M. Crampes, "Survey on social community detection," in *Social media retrieval*, Springer, 2012, pp. 65–85.

[11]  F. D. Malliaros and M. Vazirgiannis, "Clustering and community detection in directed networks: A survey," *Physics reports*, vol. 533, no. 4, pp. 95–142, 2013.

[12]  L. Despalatović, T. Vojković, and D. Vukicević, "Community structure in networks: Girvan-newman algorithm improvement," in *2014 37th international convention on information and communication technology, electronics and microelectronics (MIPRO)*, IEEE, 2014, pp. 997–1002.

[13]  J. Maluck and R. V. Donner, "A network of networks perspective on global trade," *PloS one*, vol. 10, no. 7, e0133310, 2015.

[14] P. Bedi and C. Sharma, "Community detection in social networks," *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, vol. 6, no. 3, pp. 115–135, 2016.

[15] S. Dwivedi, P. Kasliwal, and S. Soni, "Comprehensive study of data analytics tools (rapidminer, weka, r tool, knime)," in *2016 Symposium on Colossal Data Analysis and Networking (CDAN)*, IEEE, 2016, pp. 1–8.

[16] S. Fortunato and D. Hric, "Community detection in networks: A user guide," *Physics reports*, vol. 659, pp. 1–44, 2016.

[17] D. Miller, J. Sinanan, X. Wang, *et al.*, *How the world changed social media.* UCL press, 2016.

[18] S. Leavy, "Gender bias in artificial intelligence: The need for diversity and gender theory in machine learning," in *Proceedings of the 1st international workshop on gender equality in software engineering*, 2018, pp. 14–16.

[19] N. Suzor, M. Dragiewicz, B. Harris, R. Gillett, J. Burgess, and T. Van Geelen, "Human rights by design: The responsibilities of social media platforms to address gender-based violence online," *Policy & Internet*, vol. 11, no. 1, pp. 84–103, 2019.

[20] K. Varsha and K. K. Patil, "An overview of community detection algorithms in social networks," in *2020 International Conference on Inventive Computation Technologies (ICICT)*, IEEE, 2020, pp. 121–126.

[21] K. Varsha and K. K. Patil, "An overview of community detection algorithms in social networks," in *2020 International Conference on Inventive Computation Technologies (ICICT)*, IEEE, 2020, pp. 121–126.

[22] P. Vispute, "Performance evaluation of community detection algorithms in social networks analysis," *Bioscience Biotechnology Research Communications*, vol. 13, no. 14, pp. 388–393, 2020. DOI: 10.21786/bbrc/13.14/90.

[23] S. Ali, M. H. Shakeel, I. Khan, S. Faizullah, and M. A. Khan, "Predicting attributes of nodes using network structure," *ACM Transactions on Intelligent Systems and Technology (TIST)*, vol. 12, no. 2, pp. 1–23, 2021.

[24] L. Jamison, A. P. Christensen, and H. Golino, "Optimizing walktrap's community detection in networks using the total entropy fit index," 2021.

[25] S. Pan, C.-c. Yang, J.-Y. Tsai, and C. Dong, "Experience of and worry about discrimination, social media use, and depression among asians in the united states during the covid-19 pandemic: Cross-sectional survey study," *Journal of Medical Internet Research*, vol. 23, no. 9, e29024, 2021.

[26] K. J. Wong, *7 evaluation metrics for clustering algorithms*, Dec. 2022. [Online]. Available: https://towardsdatascience.com/7-evaluation-metrics-for-clustering-algorithms-bdc537ff54d2.

[27] *King University Online*, Apr. 2023. [Online]. Available: https://online.king.edu/news/psychology-of-social-media/.

[28] H. Buie and A. Croft, "The social media sexist content (smsc) database: A database of content and comments for research use," *Collabra: Psychology*, vol. 9, no. 1, p. 71 341, 2023.

[29]    *Troll Patrol Report*, [Online]. Available: https://decoders.amnesty.org/projects/troll-patrol/findings.

[30]    [Online]. Available: https://sites.google.com/site/bctnet/network-construction.

[31]    *NVIDIA Data Science Glossary*, [Online]. Available: https://www.nvidia.com/en-us/glossary/data-science/networkx/.

[32]    *Graph types - NetworkX 3.1 documentation*, [Online]. Available: https://networkx.org/documentation/stable/reference/classes/index.html.