

# Community Detection in GitHub User Networks Using AGM

Akshat Raj Saxena  
Roll Number: 2022054  
GitHub Project Link

## 1 Introduction

In this project, I analyze GitHub user networks to identify communities using the AGM algorithm. This project demonstrates the practical application of graph theory and community detection algorithms in social network analysis with a large dataset.

## 2 Dataset Description

The dataset used in this project consists of GitHub user interactions represented as edges in a graph. The data is stored in a CSV format with the following structure:

---

1	id_1,id_2
2	1,2
3	2,3
4	...

---

where each row represents an edge between two users identified by their unique IDs.

## 3 Methodology

### 3.1 Graph Creation

The first step involves creating a network graph from the edge data. We use NetworkX to construct an undirected graph where nodes represent users and edges represent connections between them. Below is the sample code.

---

```
1 import networkx as nx
2 import csv
3
4 edges_list = []
5 with open("musae_git_edges.csv", 'r') as file:
6     csv_reader = csv.DictReader(file)
7     for row in csv_reader:
8         id_1 = int(row['id_1'])
9         id_2 = int(row['id_2'])
10        edges_list.append((id_1, id_2))
11
12 G = nx.Graph()
13 G.add_edges_from(edges_list)
```

---

### 3.2 Community Detection Algorithm

We implement the Affiliation Graph Model algorithm using NetworkX's built-in functions. The algorithm works by:

1. Initially treating each node as a separate community
2. Iteratively merging communities to maximize modularity
3. Stopping when no further improvement in modularity is possible

---

```
1 from networkx.algorithms.community import
   greedy_modularity_communities
2 communities = list(greedy_modularity_communities(G))
3 modularity = nx.algorithms.community.quality.modularity(G,
   communities)
```

---

### 3.3 Community Assignment

After detecting communities, I assign each node to its respective community and create a mapping for analysis:

---

```
1 from collections import defaultdict
2
3 communities_dict = defaultdict(list)
4 for i, community in enumerate(communities):
5     for node in community:
6         communities_dict[node] = i
7 from pyspark.sql import SparkSession
8 from pyspark.sql.types import StructType, StructField, IntegerType
9
10 spark =
11     SparkSession.builder.appName("CommunityDetectionAGM").getOrCreate()
12 edges_schema = StructType([
13     StructField("id_1", IntegerType(), True),
14     StructField("id_2", IntegerType(), True)
15 ])
16 edges_df = spark.read.csv("musae_git_edges.csv",
17                             header=True,
18                             schema=edges_schema)
```

---

## 4 Results and Analysis

### 4.1 Community Structure

The algorithm successfully identified distinct communities within the network. Key metrics include:

- Number of communities detected
- Size distribution of communities
- Modularity score of the partition

## 4.2 Modularity Analysis

The modularity score provides a measure of the quality of the community detection:

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

where:

- $m$  is the number of edges
- $A_{ij}$  is the adjacency matrix
- $k_i$  is the degree of node  $i$
- $c_i$  is the community of node  $i$
- $\delta$  is the Kronecker delta function

## 5 Working with the Project

To reproduce the results:

1. Clone the repository
2. Install dependencies:

---

```
1 pip install networkx numpy pandas pyspark
```

---

3. Run the community detection script:

---

```
1 python community_detection.py
```

---

## 6 Conclusion

This project demonstrates the application of the AGM algorithm for community detection in GitHub user networks. The implementation provides both PySpark and Python solutions, making it suitable for various use cases and dataset sizes.