# Project Final Report

# Detecting Tight Communities in Facebook

*Naveenraj Palanisamy (nxp154130)*

*Manoj Kumar Natha (mxn151930)*

*Anandan Sundar (axs156730)*

*Abhishek Thangudu (axt153530)*

*Venkata Sai Praneeth Palakurthi (vxp150830)*

# Table of Contents

# 1. Introduction

The identification of tight communities in a network is very important to understand the structure of the Network. This is helpful in many areas like Medical science, Defense teams. In Medical Science this analysis can be used for detecting Cancer causing genes, in Defense for terrorist community detection. We can find the like-mindedness among the people in Social Networks.

# 2. Project Description and Algorithm Approach

## 2.1. Dataset Description

This project deals with the finding of tightness among communities across Facebook users. A large dataset has been taken from SNAP Website which consists of details about Facebook users and the relationships between them. The original IDs are replaced with new values. All the feature names are also replaced by generic names. The dataset also contain groups to signify common interests among the Users. The Users in the network are considered as Vertices of the Graph and the relationship among the Users are considered as Edges of the Graph.

## 2.2. Dataset Statistics

Dataset Link: http://snap.stanford.edu/data/egonets-Facebook.html

The following table describes the SNAP Facebook dataset as in the Dataset description:

| Dataset statistics | |
| --- | --- |
| Nodes | 4039 |
| Edges | 88234 |
| Nodes in largest WCC | 4039 (1.000) |
| Edges in largest WCC | 88234 (1.000) |
| Nodes in largest SCC | 4039 (1.000) |
| Edges in largest SCC | 88234 (1.000) |
| Average clustering coefficient | 0.6055 |
| Number of triangles | 1612010 |
| Fraction of closed triangles | 0.2647 |
| Diameter (longest shortest path) | 8 |
| 90-percentile effective diameter | 4.7 |

The Facebook dataset resembles an Undirected Graph. This is because two Users in Facebook are mutually connected if they are friends unlike in Twitter or Google+ where a User follows another User which signifies a directed graph. There is no difference between weakly and strongly connected components in an undirected graph. This is the reason the sub graph of largest WCC and SCC shows same number of nodes in the above table. The Average clustering coefficient of the whole signifies the overall tightness of the network. As all nodes in the SCC are equal to the total number of nodes in the graph, there is only one connected component in the given data. Since it is an undirected graph, SCC and WCC are same. So, detecting strongly connected components in the graph is same as detecting connected components in this graph. Also, since we have only one connected component, we are increasing our project's scope to finding cliques in the graph. A completely connected sub-graph in the given network is called a clique.

The Clique (Completely connected sub-component) is different from a strongly connected sub-component in the following way. In a completely connected sub-component, there is an edge between all the distinct nodes in the sub-component whereas in strongly connected sub-component, a path between any two distinct nodes in the sub-component exists.

## 2.3. Project Approach

All the analysis requires constructing and analyzing a network which can be modeled as a Graph. Apache Spark has provided an API for graphs and parallel computation in graphs which is GraphX. GraphX API is used in this project for constructing the graph and for analyzing the graph. In addition to this, JGrapht API is used. JGrapht has additional functionalities for graph processing. Here we use JGrapht for finding the cliques across the graph.

As a first step, a Graph is built using the Facebook dataset with Users as the vertices and the relationship among them as edges. The graph is built in such a way that if there is an edge between two vertices, it means the Users are friends in Facebook. All the edge weights are unique and taken as 1.

The second step is finding the strongly connected components and cliques (completely connected sub-graph) in the network-graph. To find this, we are using Bron-Kerbosch algorithm which is in Jgrapht. We get all the cliques in the graph listed. Now, we find the clustering coefficient of all vertices in the graph using the triangle count. Clustering coefficient of a vertex is measured based on number of triangle counts formed for each vertex to total number of triangles that can be formed at a vertex ($^{n}C_2$ where n is

number of neighbors for that Vertex). Then, tightness of each clique is determined using the clustering coefficient. Clustering coefficient of a clique is the average of the clustering coefficients of each vertex present in the clique. The value of clustering coefficient ranges in between 0 and 1. Based on these values we can determine how tight the given clique is.

The third step is determining the tightness for each clique in the overall network and ordering the cliques based on the tightness. These results are recorded and reported in next section.

## 3. Experimental Analysis and Results

DFS (Depth First Search) traversal is used for finding the strongly connected components in a given graph network. GraphX have provided a direct method for finding the strongly connected components. In this dataset, entire graph is a strongly connected component.

Main logic of code:

Output showing the average clustering coefficient of the whole graph:



Bron-Kerbosch algorithm can be used for finding maximal cliques in an undirected graph. JGraphT has provided direct method for finding cliques in a graph.

Main logic of code:

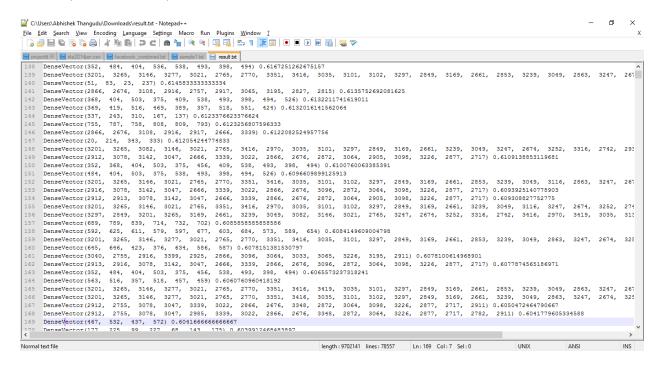Output showing the cliques in the decreasing order of co-efficient of tightness:

Output saved in text files:

i)    All the cliques in the graph (file: sample1.txt) :



ii)    Final result of cliques in the graph with decreasing order of co-efficient of tightness (File : result.txt):

## 4. Future work

The analysis can further be extended by comparing the features of each User in a given tight community and thereby determining which feature made them to be in a tight community. We can also measure connectivity of each Facebook group by considering only sub-graph of a group separately and measuring the tightness. But in the dataset provided, groups of different ego networks are not connected properly; they are just named in a sequential manner. There is no way connecting the members in one ego network to other.

## 5. Conclusion

Finding Cliques is one of the best ways in determining the tight communities in a given network. Strongly connected components can just tell if there is a path between any two distinct users that can be connected whereas in a Clique there is a direct connection among every distinct User in the sub-network.

## 6. References

1. http://spark.apache.org/graphx/
2. http://jgrapht.org/javadoc/org/jgrapht/alg/BronKerboschCliqueFinder.html
3. http://stackoverflow.com/questions/31217642/finding-cliques-or-strongly-connected-components-in-apache-spark-using-graphx
4. https://en.wikipedia.org/wiki/Bron%E2%80%93Kerbosch_algorithm