

Web and Social Computing (IT 752)

Submitted by: Sampat Kr Ghosh

Roll Number: 192IT020

Submitted to: Dr. Sowmya Kamath S.

Lab Assignment 3

Dataset 1: Facebook

Dataset 2: Gnutella P2P Network, August 4, 2002

Dataset 3: Twitter

Part 1

The datasets mentioned above were used to measure the centrality of nodes - Degree centrality, Closeness centrality, Betweenness centrality, and Eigenvector centrality. For plotting the distributions, normal and logarithmic scale both are used as for some distributions that are in normal scale, I was not getting a good plot.

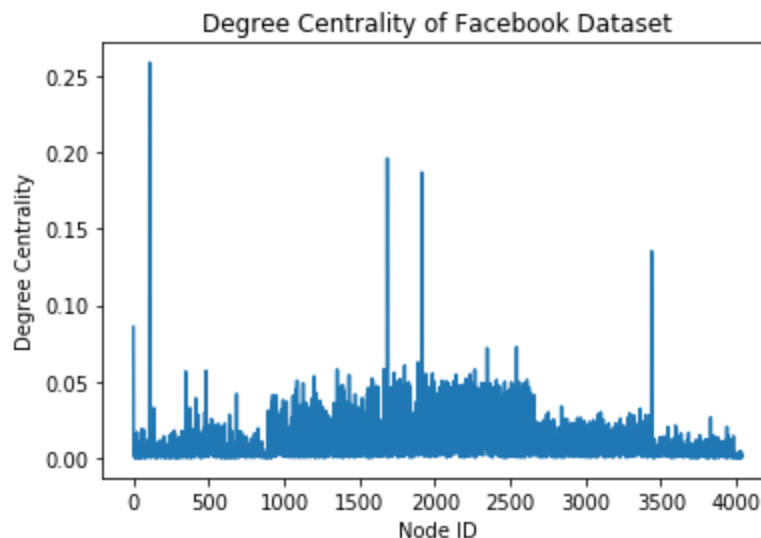
Degree Centrality

Degree centrality is the simplest centrality measure to compute. Recall that a node's degree is simply a count of how many social connections (i.e., edges) it has. The degree centrality for a node is simply its degree. For degree centrality, higher values mean that the node is more central. Degree centrality shows how many connections a person has. They may be connected to lots of people at the heart of the network, but they might also be far off on the edge of the network.

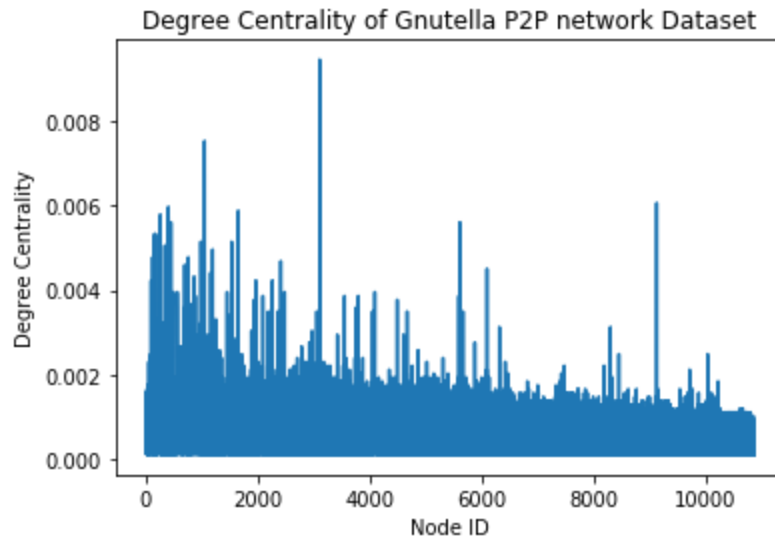
Score (C_d) = number of edges attached to the node.

Standardized Score = $C_d / (n-1)$, where n is the total number of nodes.

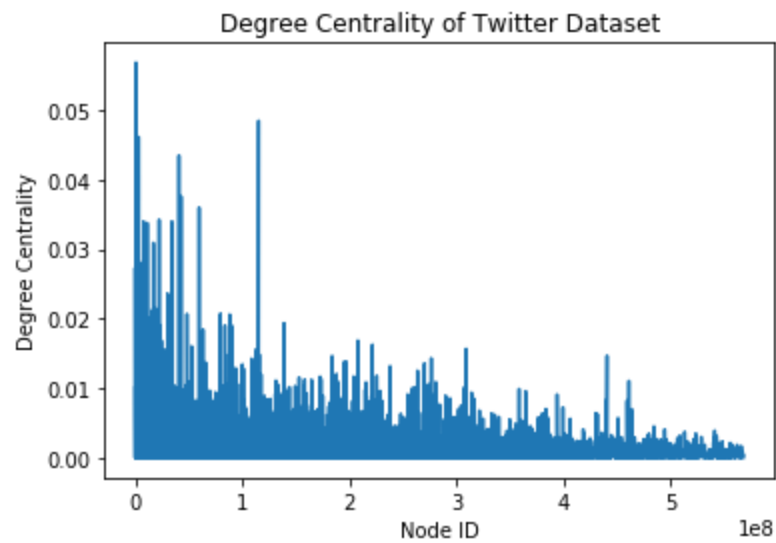
Degree Centrality V/S Node ID plots are given below:



Facebook

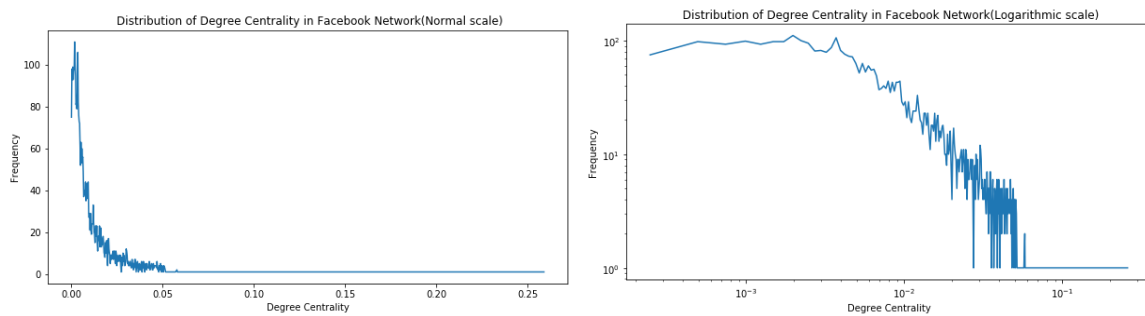


Gnutella P2P Network

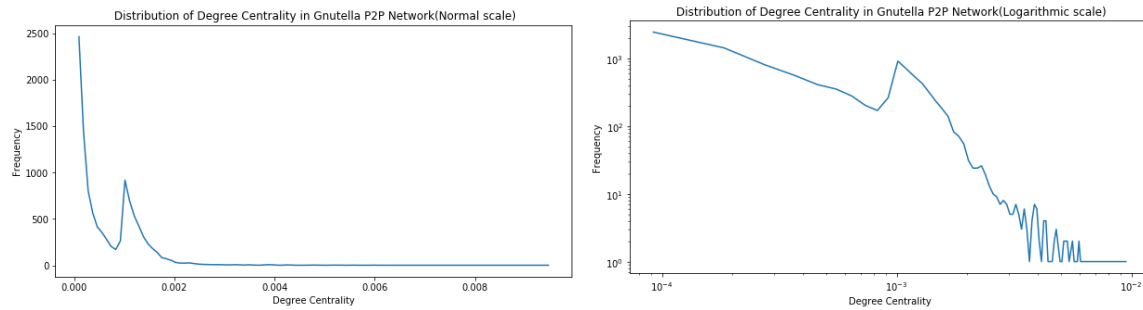


Twitter

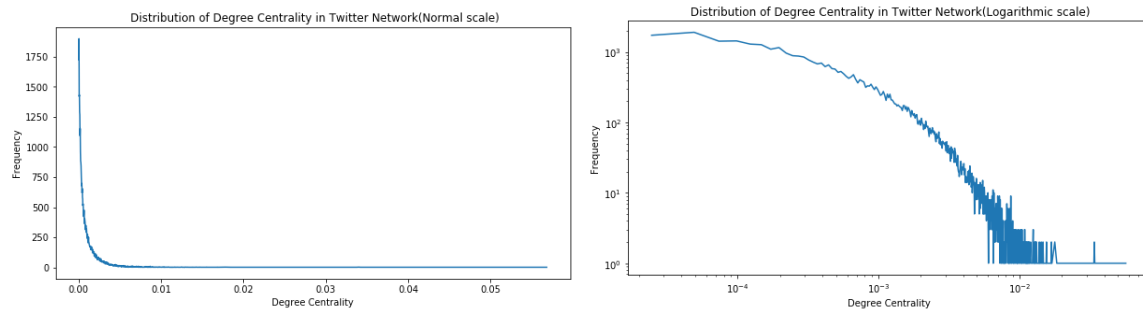
Degree centrality distribution for all the datasets are plotted below:



Facebook



Gnutella P2P Network



Twitter

Observation

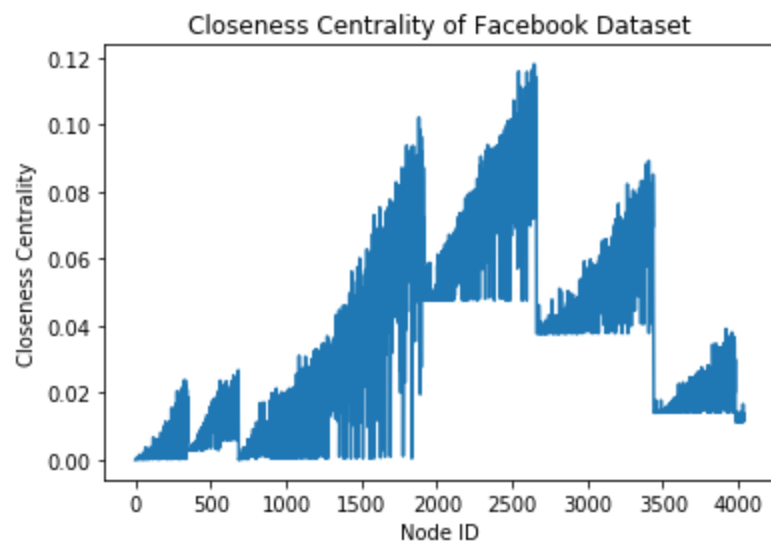
Similar to Degree distribution the degree centrality distribution also shows power law of distribution ie. very few nodes have higher degree centrality and the majority of nodes have very small degree centrality.

Closeness Centrality

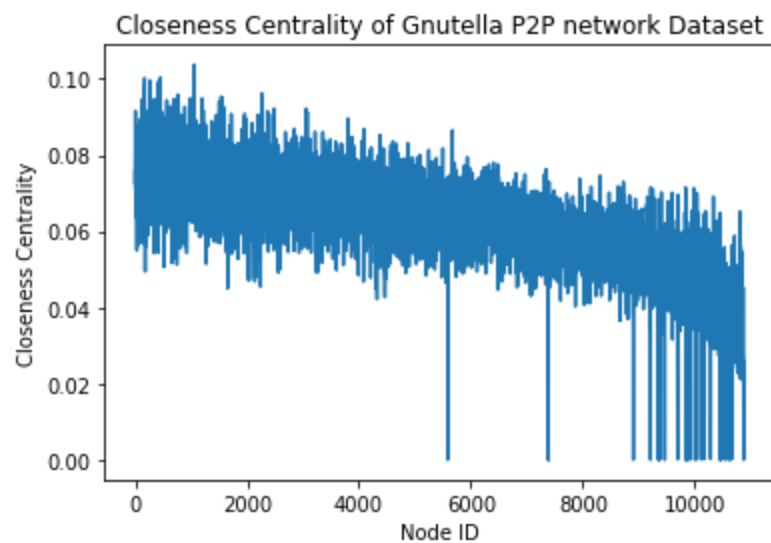
Closeness centrality indicates how close a node is to all other nodes in the network. It is calculated as the average of the shortest path length from the node to every other node in the network. It is defined as the reciprocal of the average shortest path length.

Score = $1 / \text{avg} (L(n,m))$, where $L(n,m)$ is the length of the shortest path between two nodes n and m .

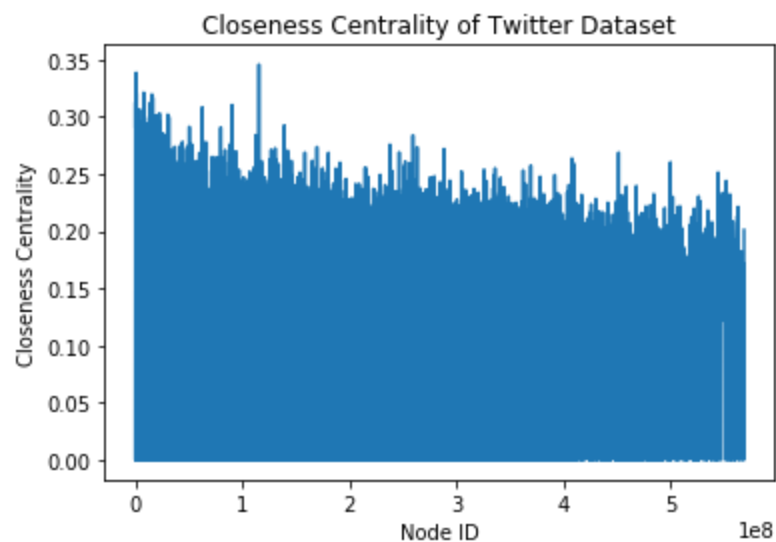
Closeness centrality V/S Node ID plots are given below:



Facebook

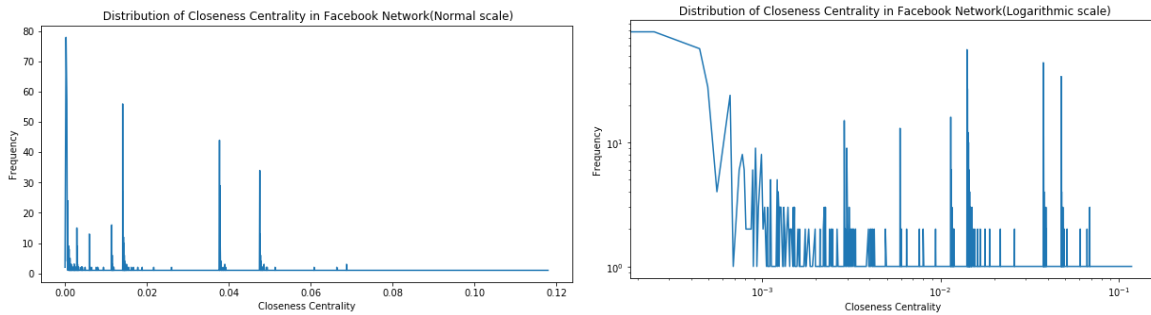


Gnutella P2P Network

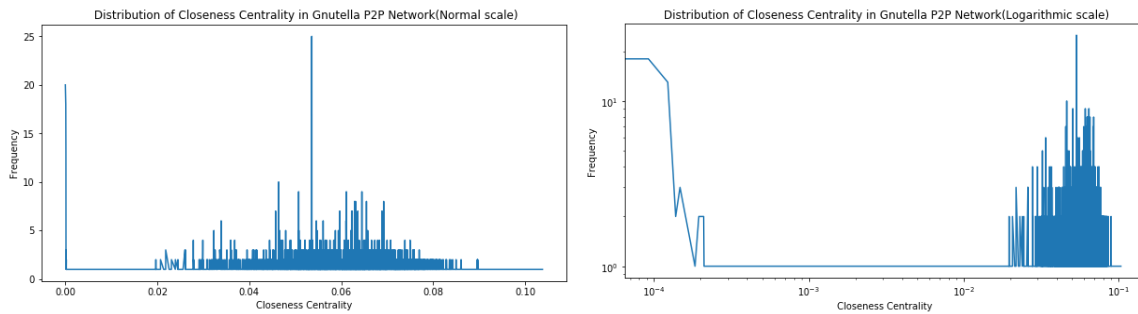


Twitter

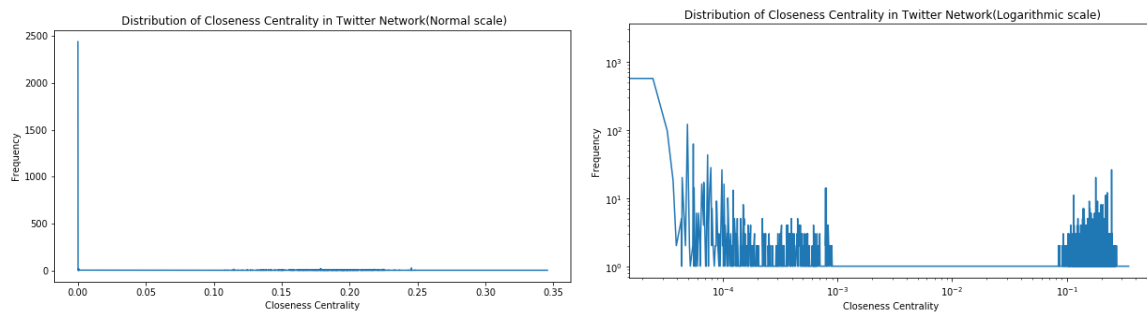
Closeness centrality distribution for all the datasets are plotted below:



Facebook



Gnutella P2P Network



Twitter

Observation

Closeness centrality can help find good 'broadcasters', but in a highly-connected network, you will often find all nodes have a similar score. What may be more useful is using Closeness to find influencers in a single cluster. As the above datasets are not highly connected, we have dissimilar closeness centrality score.

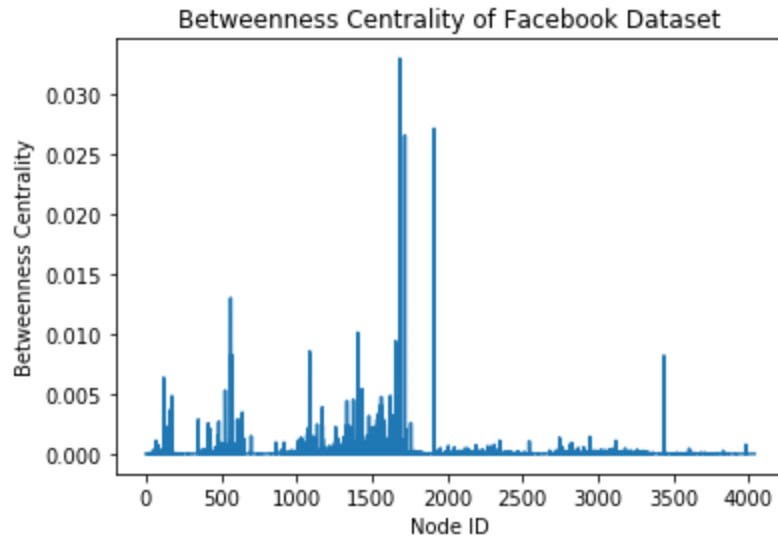
Betweenness Centrality

Betweenness centrality measures how important a node is to the shortest paths through the network. To compute betweenness for a node N, we select a pair of nodes and find all the shortest paths between those nodes. Then we compute the fraction of those shortest paths that include node N.

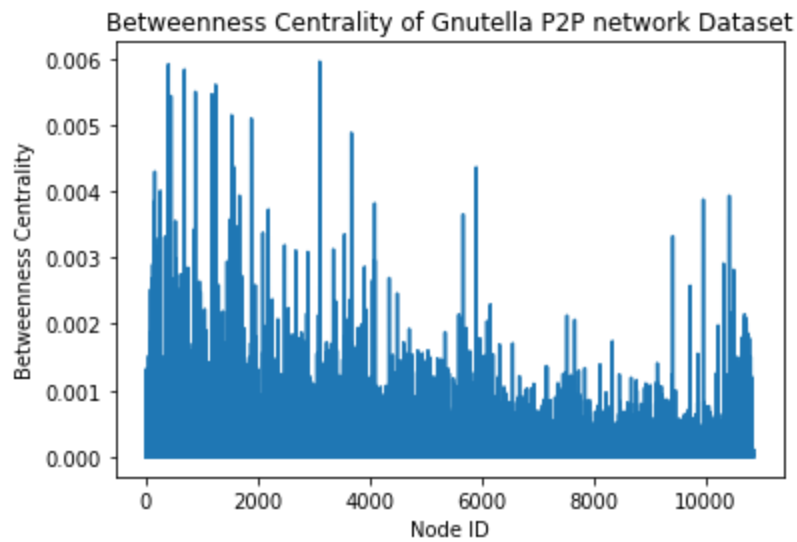
Score (C_b) = $\Sigma \sigma_{st}(v) / \sigma_{st}$, where $\sigma_{st}(v)$ is the number of shortest paths from s to t that v lies on and σ_{st} is the number of shortest paths from s to t.

Standardized Score = $C_b / [(n-1)(n-2)/2]$

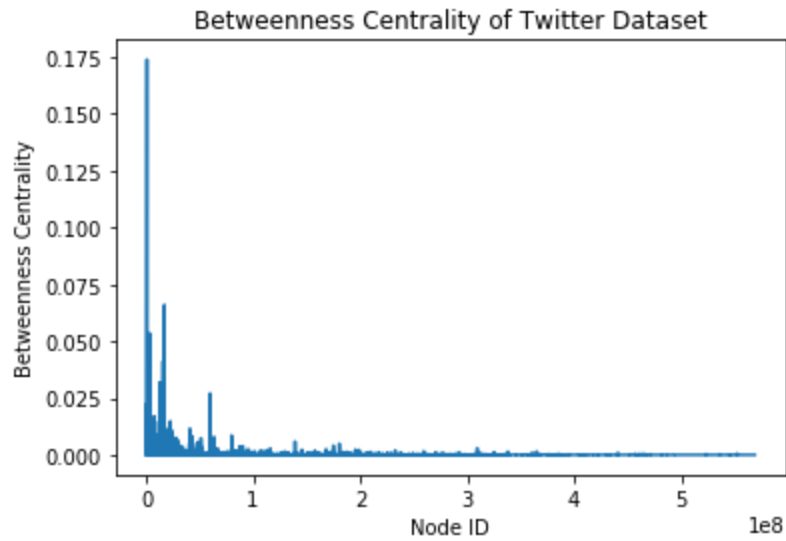
Betweenness centrality V/S Node ID plots are given below:



Facebook

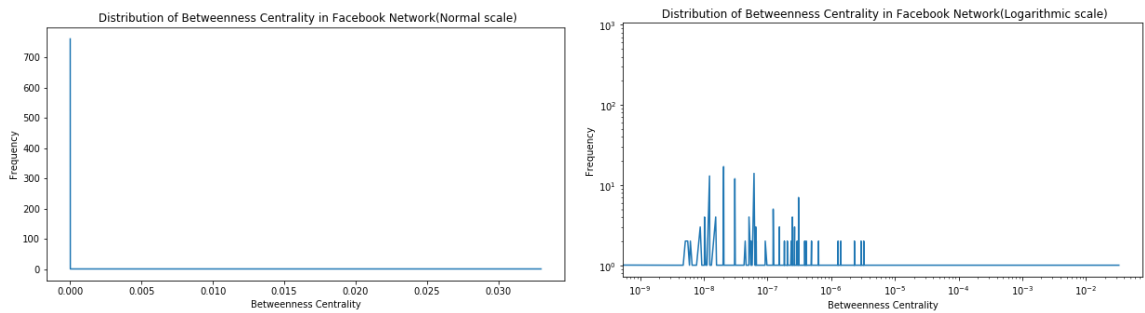


Gnutella P2P Network

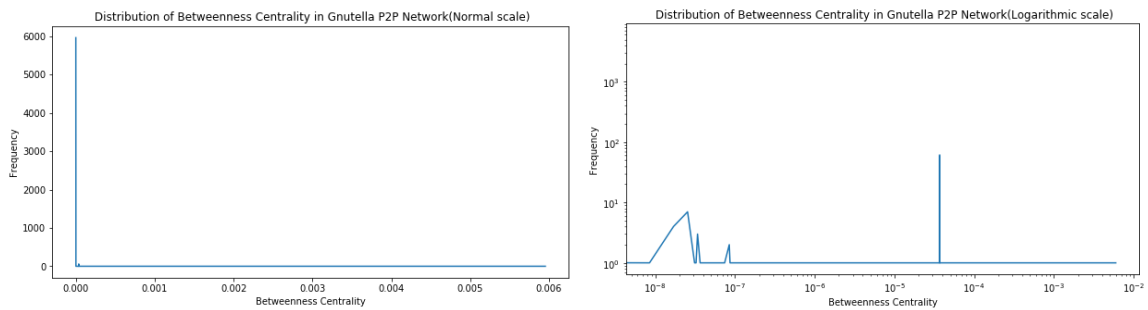


Twitter

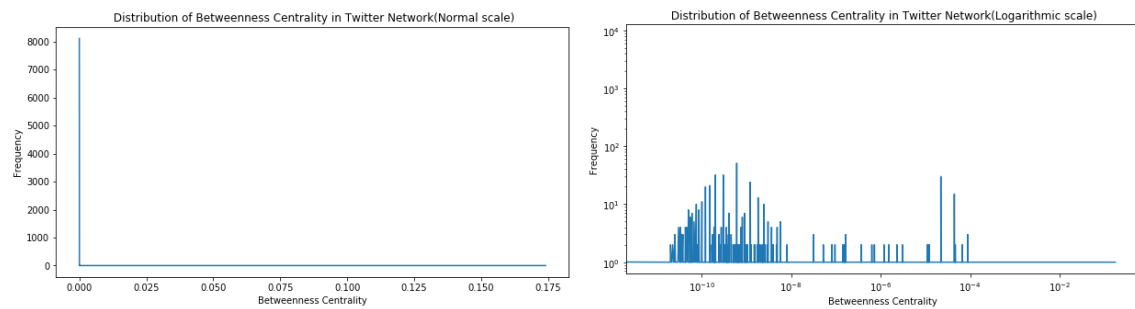
Betweenness centrality distribution for all the datasets are plotted below:



Facebook



Gnutella P2P Network



Twitter

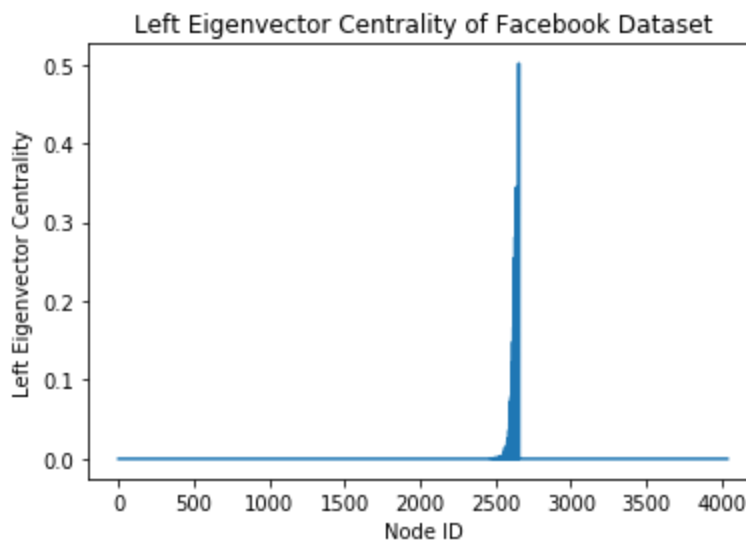
Observation

A high betweenness centrality could indicate someone holds authority over disparate clusters in a network, or just that they are on the periphery of both clusters. As we can see from the above plots, there are few nodes that have high betweenness centrality. From this we can conclude that these individuals influence the flow around a system.

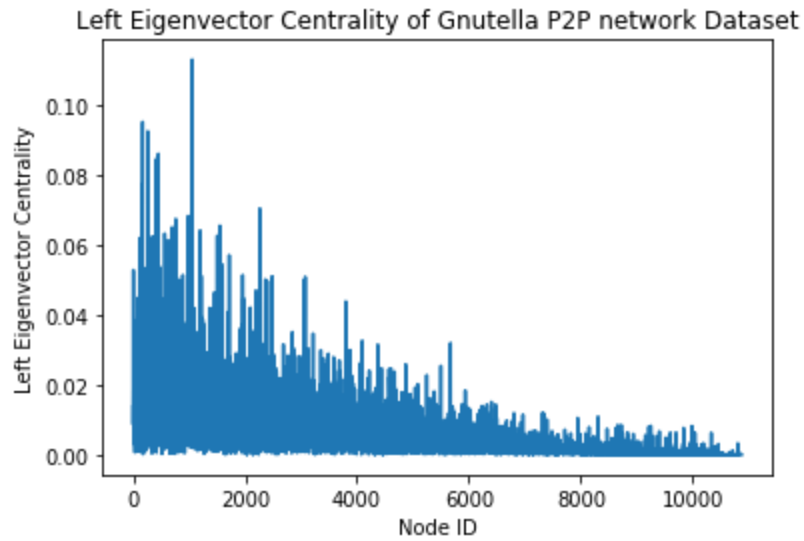
Eigenvector Centrality

Eigenvector centrality is a measure of the influence of a node in a network. It assigns relative scores to all nodes in the network based on the concept that connections to high-scoring nodes contribute more to the score of the node in question than equal connections to low-scoring nodes. Google's PageRank and the Katz centrality are variants of the eigenvector centrality. Since the datasets used here are directed graphs, we need to compute Left Eigenvector (Incoming edge) Centrality and Right Eigenvector Centrality (outgoing edges).

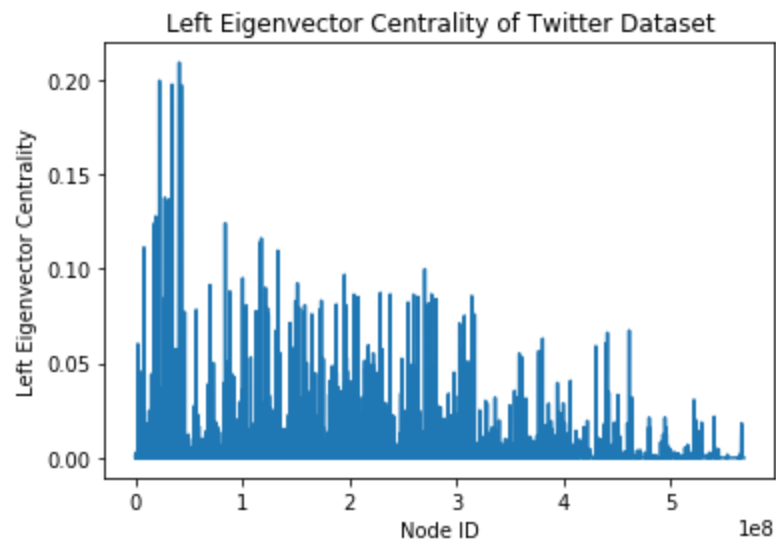
Left eigenvector centrality V/S Node ID plots are given below:



Facebook

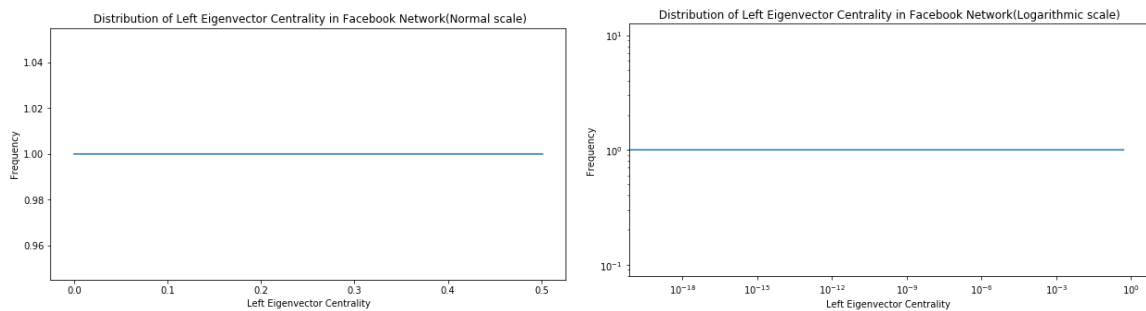


Gnutella P2P Network

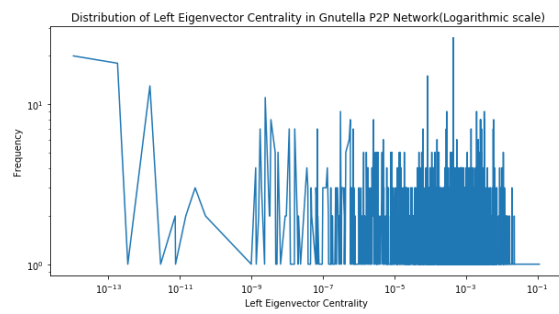
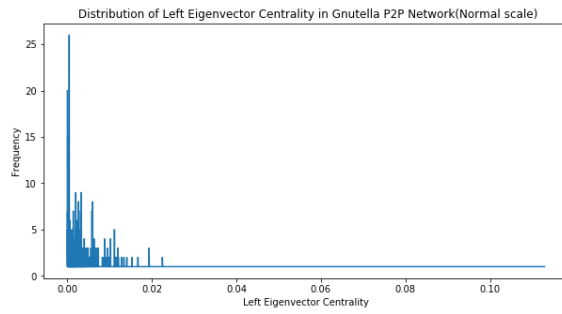


Twitter

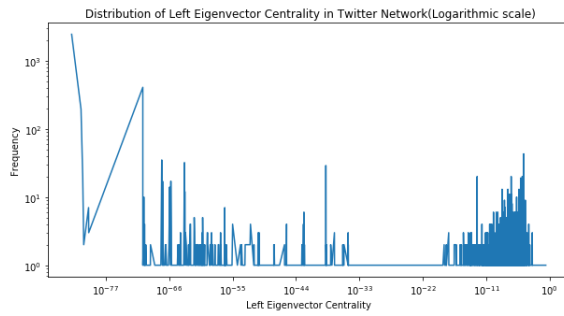
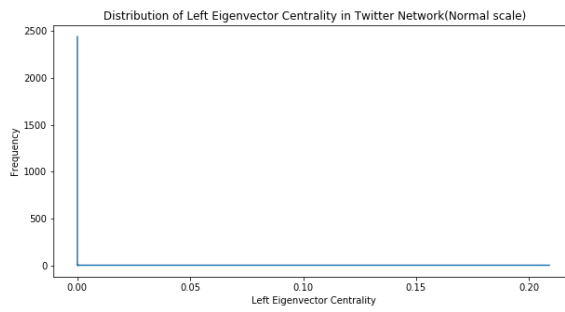
Left eigenvector centrality distribution for all the datasets are plotted below:



Facebook

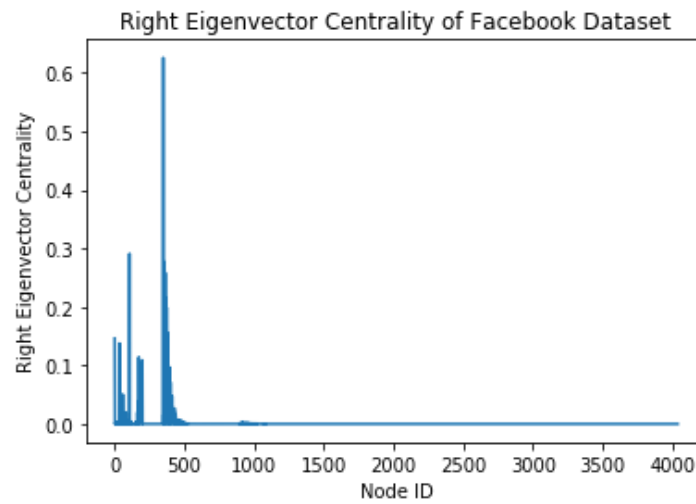


Gnutella P2P Network

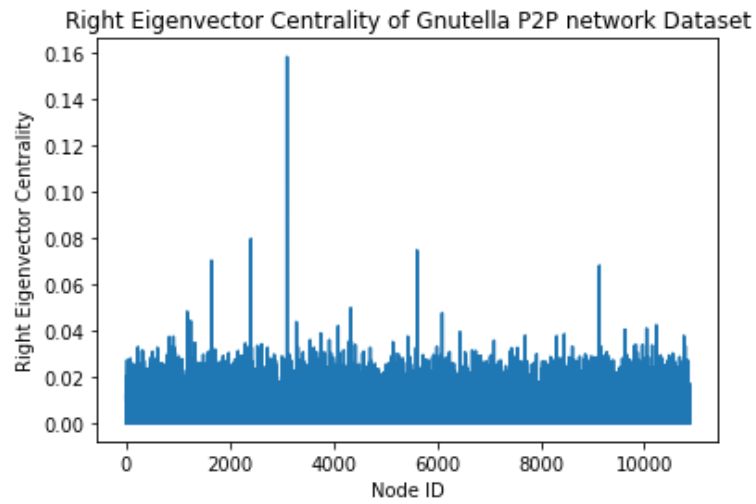


Twitter

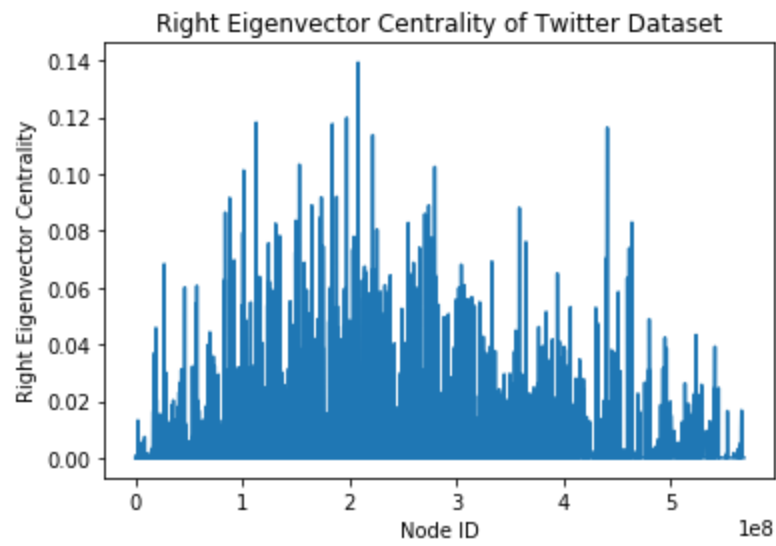
Right eigenvector centrality V/S Node ID plots are given below:



Facebook

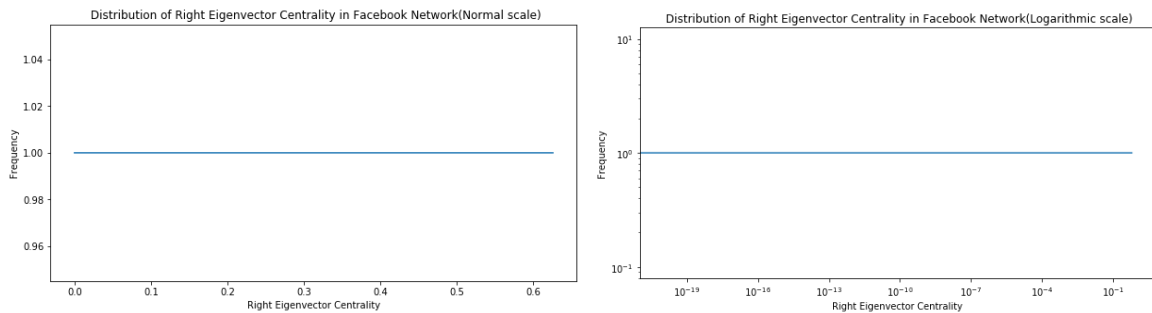


Gnutella P2P Network

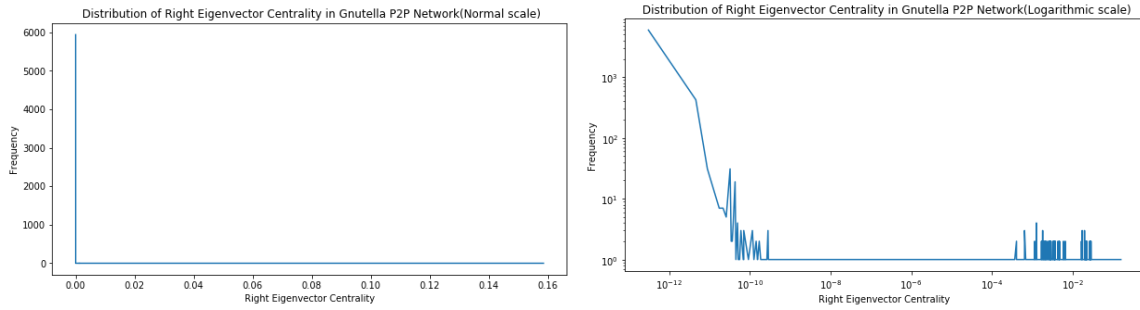


Twitter

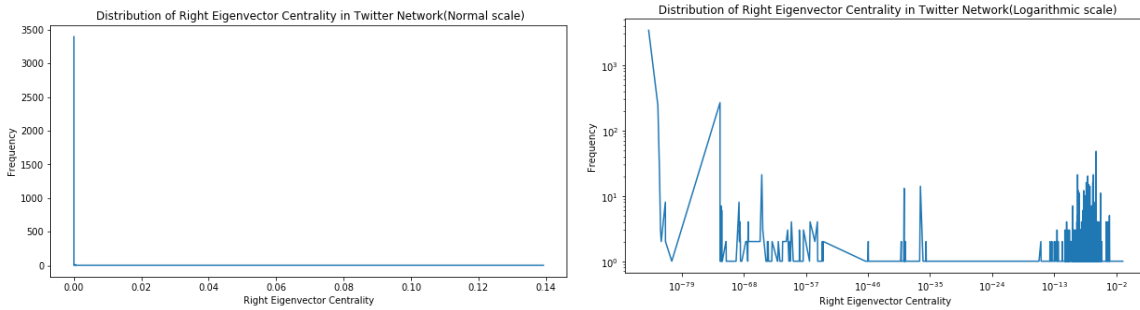
Right eigenvector centrality distribution for all the datasets are plotted below:



Facebook



Gnutella P2P Network



Twitter

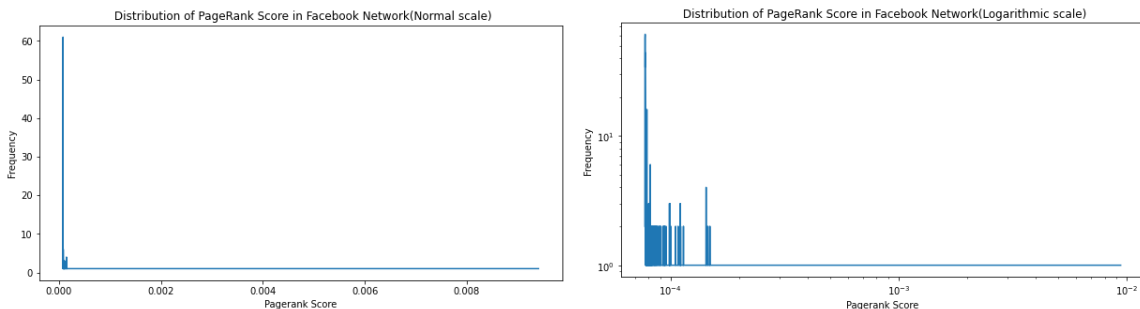
Observation

By calculating the extended connections of a node, Eigenvector Centrality can identify nodes with influence over the whole network, not just those directly connected to it. As we know that eigenvector centrality is a measure of the influence of a node in a network, from the above plots we can observe that very few nodes have high centrality and large numbers of nodes have low centrality. To interpret social networking websites, the nodes with high centrality are potential influencers.

Part 2

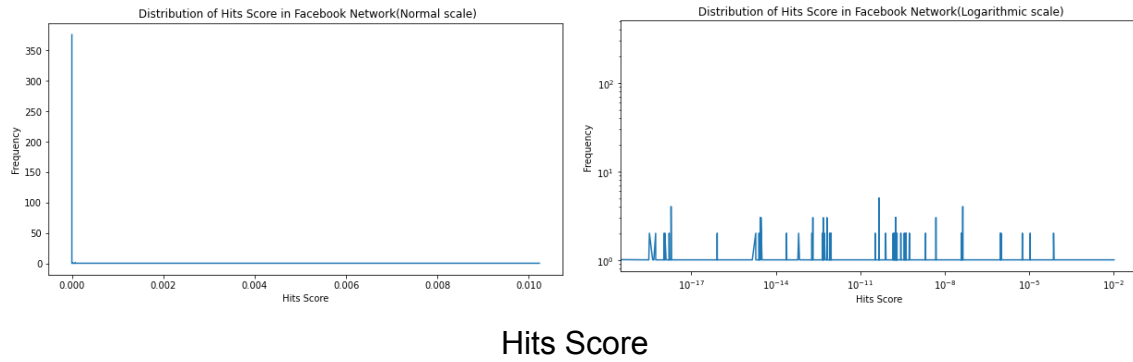
The datasets mentioned above were used to implement and compare PageRank and HITS Score algorithms. For plotting the distributions, normal and logarithmic scale both are used as for some distributions that are in normal scale, I was not getting a good plot.

Facebook Dataset

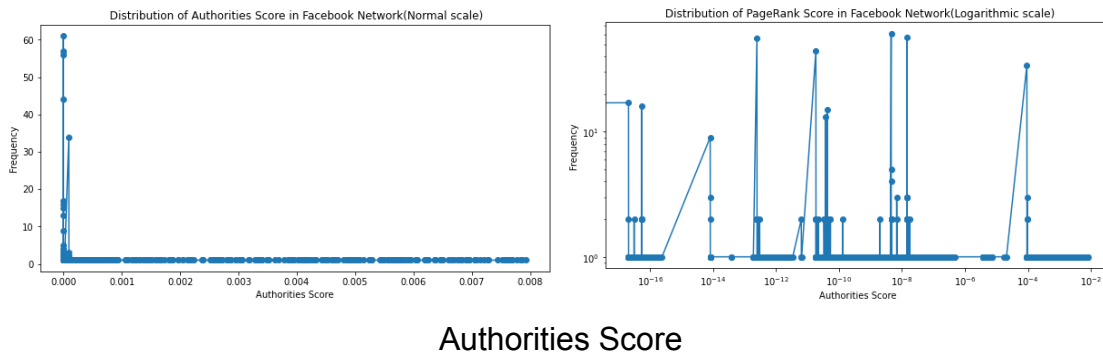


PageRank

Max PageRank Score = 0.00940916451585816
Min PageRank Score = 7.724494426592372e-05

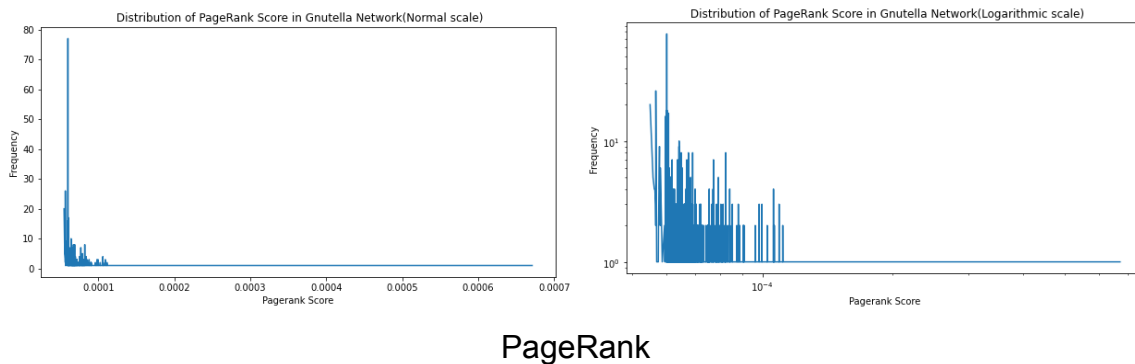


Max Hits Score = 0.010229403947448484
Min Hits Score = 0.0

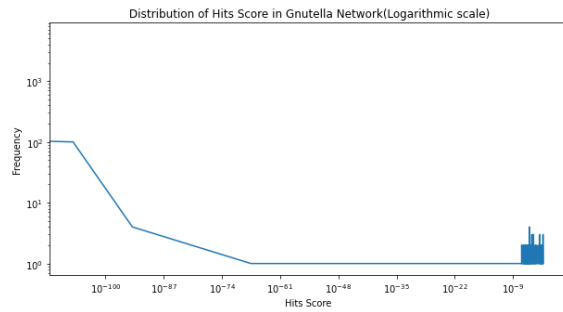
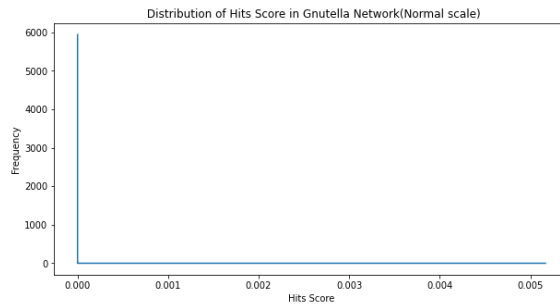


Max Authorities Score = 0.007932133892198798
Min Authorities Score = 0.0

Gnutella P2P Network



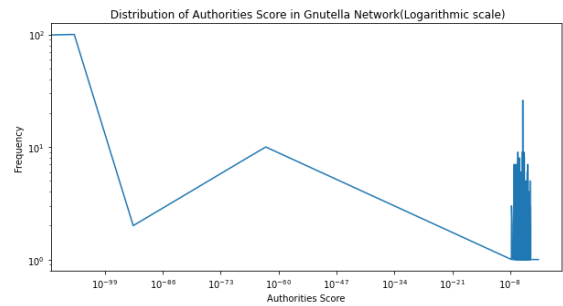
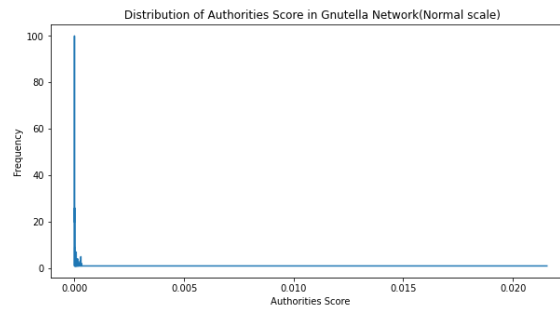
Max PageRank Score = 0.0006711727183638689
Min PageRank Score = 5.499573860784478e-05



Hits Score

Max Hits Score = 0.005167046979475104

Min Hits Score = 0.0

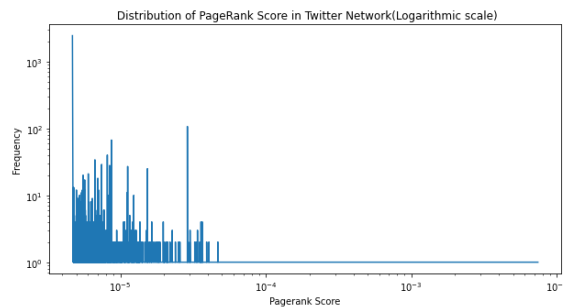
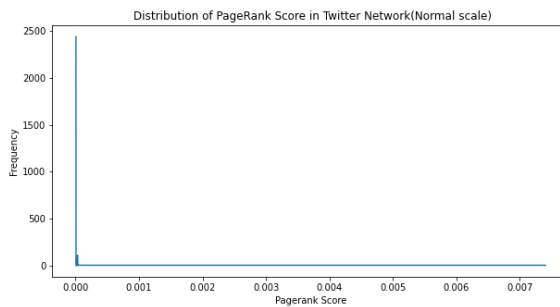


Authorities Score

Max Authorities Score = 0.021553778629464077

Min Authorities Score = 0.0

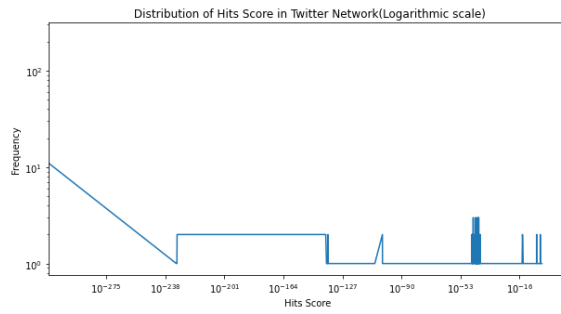
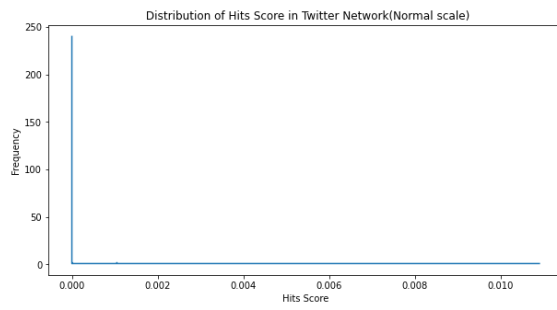
Twitter Dataset



PageRank

Max PageRank Score = 0.007408496105825751

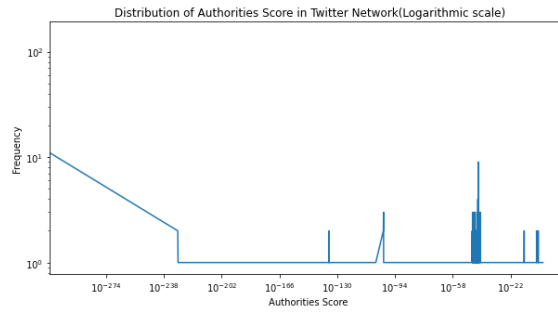
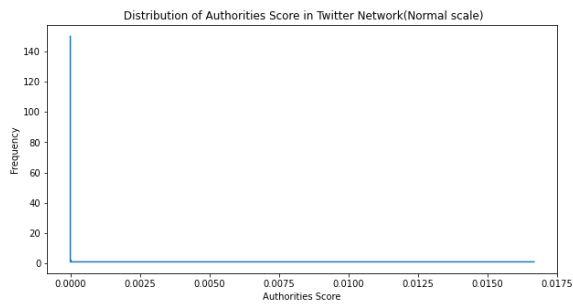
Min PageRank Score = 4.689332128677982e-06



Hits Score

Max Hits Score = 0.010907823237657264

Min Hits Score = 0.0



Authorities Score

Max Authorities Score = 0.01668608174631693

Min Authorities Score = 0.0