

CS286 – Social Network Analysis  
Term Project Report  
Ujjawal Garg (SJSU ID: 011917334)  
ujjawal.garg@sjsu.edu

# I. Introduction

Social networks enable propagation of ideas and information through interactions between individuals. One of the most important concept is to identify a set of nodes that can propagate these ideas to the largest range on the network. This is achieved by ranking the nodes based on some metric. Generally, this metric is some sort of centrality measure, that defines how *central* the node is in the network. Three most common types of centrality measures are:

1. **Degree Centrality:** Rank the nodes based on their degree
2. **Closeness Centrality:** Rank the nodes based on their distances to other nodes
3. **Betweenness Centrality:** Rank the nodes based on how many shortest paths go through them.

In order to evaluate these methods, we need some model to determine the influence. For this project, we follow the research done by Chen et. al [1]. In this paper, author used SIR (Susceptible-Infected-Recovered) model [2] for this purpose. Section II describes this model and how it can be used for evaluation. Section III introduces the **Local Centrality** measure introduced by [1]. Section IV contains the results of the analysis and conclusion of this report.

## II. SIR Model

In this model, each node can be in any of three states in it's lifetime:

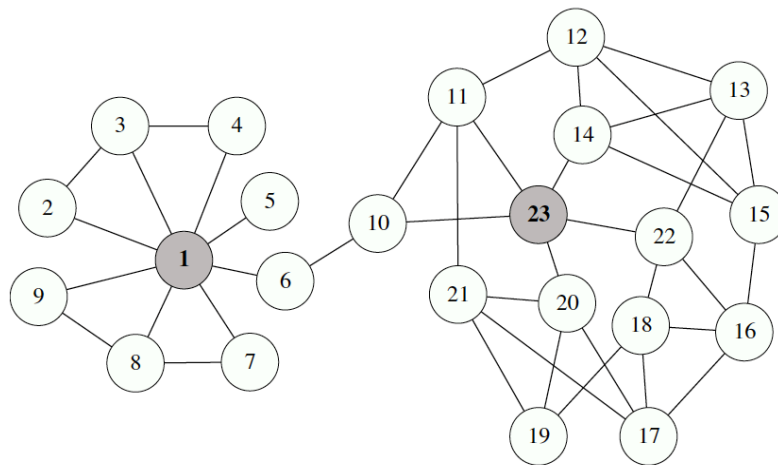
1. **Susceptible:** Not yet Infected
2. **Infected:** Will infect a randomly chosen neighbour node at each time step
3. **Recovered:** Immune to disease and will not take part in spreading infection

Initially all nodes are in Susceptible state. We select the node that we believe to be most influential and set its state to infected. Now, after each time step, for every infected node, a randomly chosen neighbour node will get infected. Infected nodes

recover with probability  $1/k$ , i.e., each node will infect an average of  $k$  nodes before becoming immune. Here  $k$  is the average degree of the network. Let  $F(t)$  denote the total number of infected and recovered nodes at time  $t$ . This value will increase monotonically and eventually become stable, once every node has been recovered.

### III. Local Centrality

Degree centrality is very simple and easiest to compute, but is of little relevance as it considers only very limited information. Consider the network show in Fig. 1.



**Fig. 1.** An example network consisted of 23 nodes and 40 edges. Although node 23 has lower degree than node 1, its influence may be even higher.

Although node 1 has highest degree centrality, the disease, if it origins at node 1, may not spread the fastest or the most broadly since all neighbours of this node have very low degree. Node 23 on the other hand has faster spreading and larger influence.

Betweenness and Closeness centrality provide much better ranking results as they consider global overview of the network. However, they have very high computation cost:  $O(n^3)$

The proposed Local centrality measure acts as a trade-off between this low-relevance and high computation cost. The local centrality  $CL(v)$  of node  $v$  is defined as:

$$Q(u) = \sum_{w \in \Gamma(u)} N(w),$$

$$C_L(v) = \sum_{u \in \Gamma(v)} Q(u),$$

where  $\Gamma(u)$  is the set of the nearest neighbours of node  $u$  and  $N(w)$  is the number of the nearest and the next nearest neighbours of node  $w$ . Since to calculate  $N(w)$  requires traversing node  $w$ 's neighbourhood within two steps, the computational complexity of local centrality is  $O(n\langle k \rangle^2)$  which grows linearly with the size of a sparse network.

## IV. Results

To analyse the performance of this new method, following datasets were used:

1. Netscience [3] : The network of co-authorships between scientists who are themselves publishing on the topic of networks. There are 379 nodes and 914 edges in this graph.
2. Email [4] : the network of e-mail interchanges between members of the University Rovira i Virgili (Tarragona)

Using the python-igraph and matplotlib library, graph was plotted between  $F(t)$  vs Centralities. Fig 2 & Fig 3 show the results for Email and Netscience dataset respectively. It is clear from these graphs that both closeness and local centralities outperform degree and betweenness metrics. While it is difficult to say which is better among closeness and local centralities, the much lower computation cost of local centrality clearly makes it the winner.

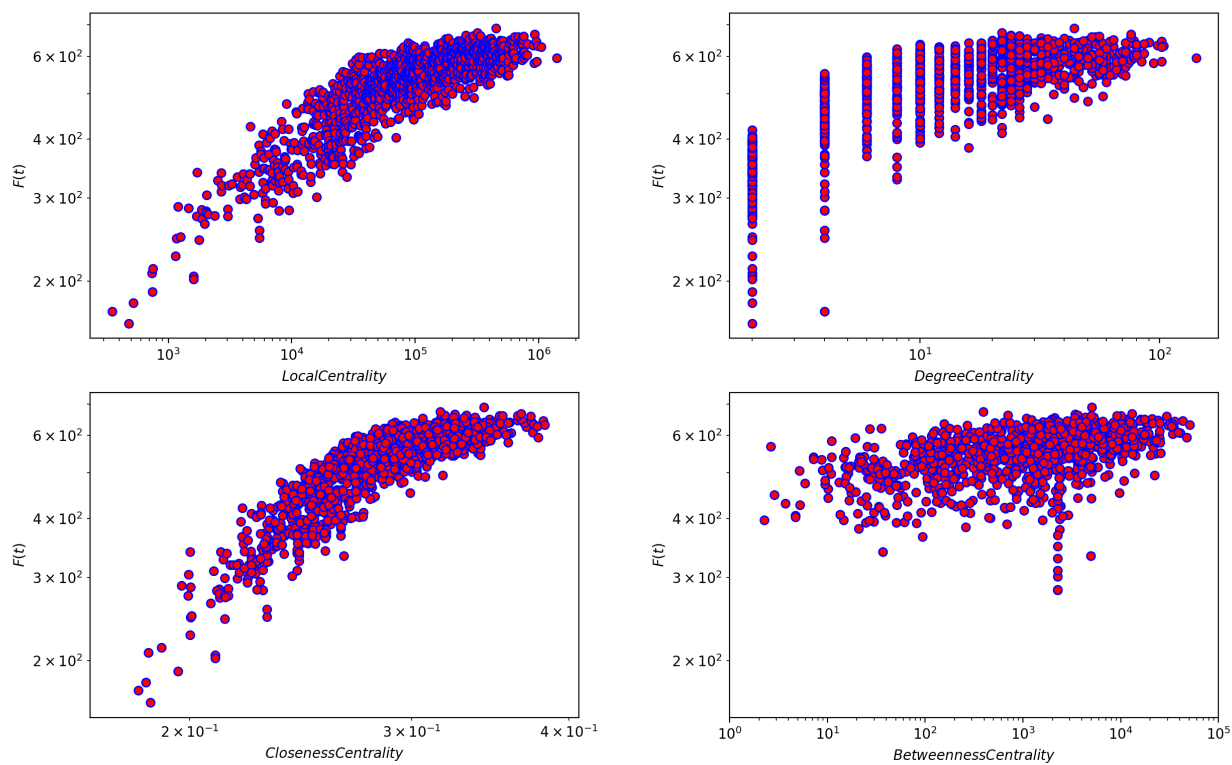


Fig 2 (Email Dataset)

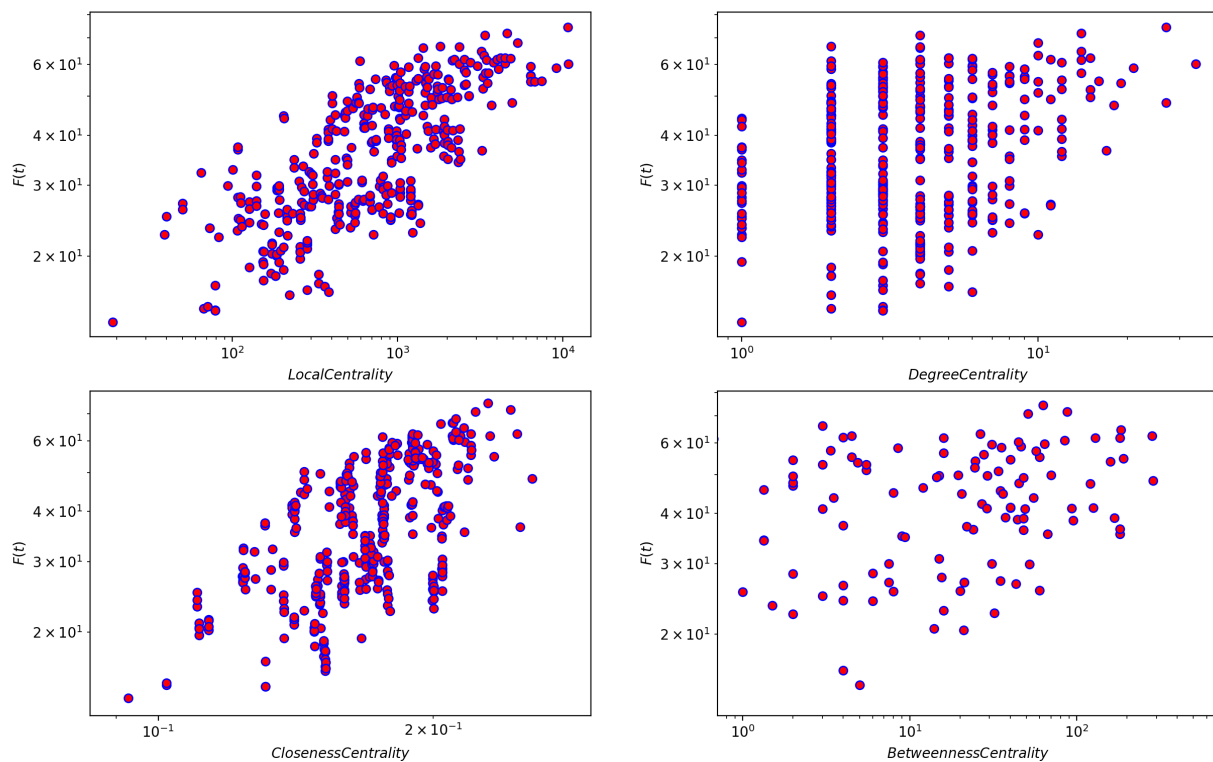


Fig 3 (Netscience Dataset)

# References

- [1] Ming-Sheng Shang Yi-Cheng Zhang Tao Zhou Duanbing Chen, Linyuan L. Identifying influential nodes in complex networks. 2012.
- [2] R.M. Anderson, R.M. May, B. Anderson, Infectious Diseases of Humans: Dynamics and Control, Oxford University Press, USA, 1992.
- [3] M. E. J. Newman, Finding community structure in networks using the eigenvectors of matrices, Preprint physics/0605087 (2006) Link: <http://nrvis.com/download/data/ca/ca-netscience.zip>
- [4] R. Guimer., L. Danon, A. Daz-Guilera, F. Giralt, A. Arenas, Self-similar community structure in a network of human interactions, Phys. Rev. E 68 (2003) 065103. Link: <http://deim.urv.cat/~aarenas/data/xarxes/email.zip>