# DATA MINING

## ASSIGNMENT-3: CLUSTERING ANALYSIS FOR COMPLEX NETWORKS

**TEAM MEMBERS:**

Srikanth Madduri Venkata
Ssmaddur
50134851

Arjun Sundaresh
arjunsun
50169917

Sridhar Vadlamani
sridharv
50168092

## Markov Clustering Algorithm

The Markov Cluster Algorithm is a fast and scalable unsupervised cluster algorithm for graphs or networks. The MCL algorithm simulates flow within a graph and promotes flow in a highly connected region and demotes otherwise, thus revealing natural groups within the graph.

In some cases, when the hub has a low curvature, many of the vertices adjacent to the hub are removed from the cluster that it represents. However, in order to make up for that, the hub cluster will absorb many other vertices—some of which are not directly connected to the hub itself—to form a large-sized core cluster. The same can be observed in ATT network as well.

Markov Clustering algorithm takes O(n3) time while finding the power or while expanding it. And O(n2) time while inflating. This is a lot of time and it can be reduced by pruning appropriately.

# Implementation

### Creating the adjacency matrix:
1. Initially we have given an index for each node between 0 and n-1 where n is the number of distinct nodes in the network.
2. As we parse through the dataset the interaction between two nodes has been marked as 1 with the help of their indices.
   Eg: if nodes "a" and "b" are interacting with each other, then if the indices of the nodes are 2 and 6 respectively, then the cells [2][6] and [6][2] of the adjacency matrix has been marked as 1.

### Adding Self loops:
3. We have also added self-loops for every node in the network in the adjacency matrix.

### Normalizing the adjacency matrix:
4. Once the self-loops are added, we have normalized the matrix. By taking the sum of the data in every column and then dividing every element in it with the sum.

### Expand the Matrix:
5. After normalizing the matrix, we have expanded the matrix by calculating it's power of 'n'. In our datasets, the 'n' value ranged from 30 to 150.

### Inflating the Matrix:

6. Once the matrix is expanded, it is inflated by calculating the power of 'x' for each and every element in the matrix.

### Normalizing the Matrix:
7. Again, after the inflation of the matrix, the matrix is normalized as described above.

**Iterating till Convergence**

8. The steps 5 to 7 are repeated till the matrix is converged.

   **Injection of results into pajek**

   9. Once the matrix is converged, the clusters are formed from the resultant matrix.
   10. A cluster is actually a row in the matrix whose members are non-zero elements in the matrix.

## Observations:

From the visualization it can be observed that as "expand" and "inflate" values are increased the number of clusters are decreased and vice-versa.
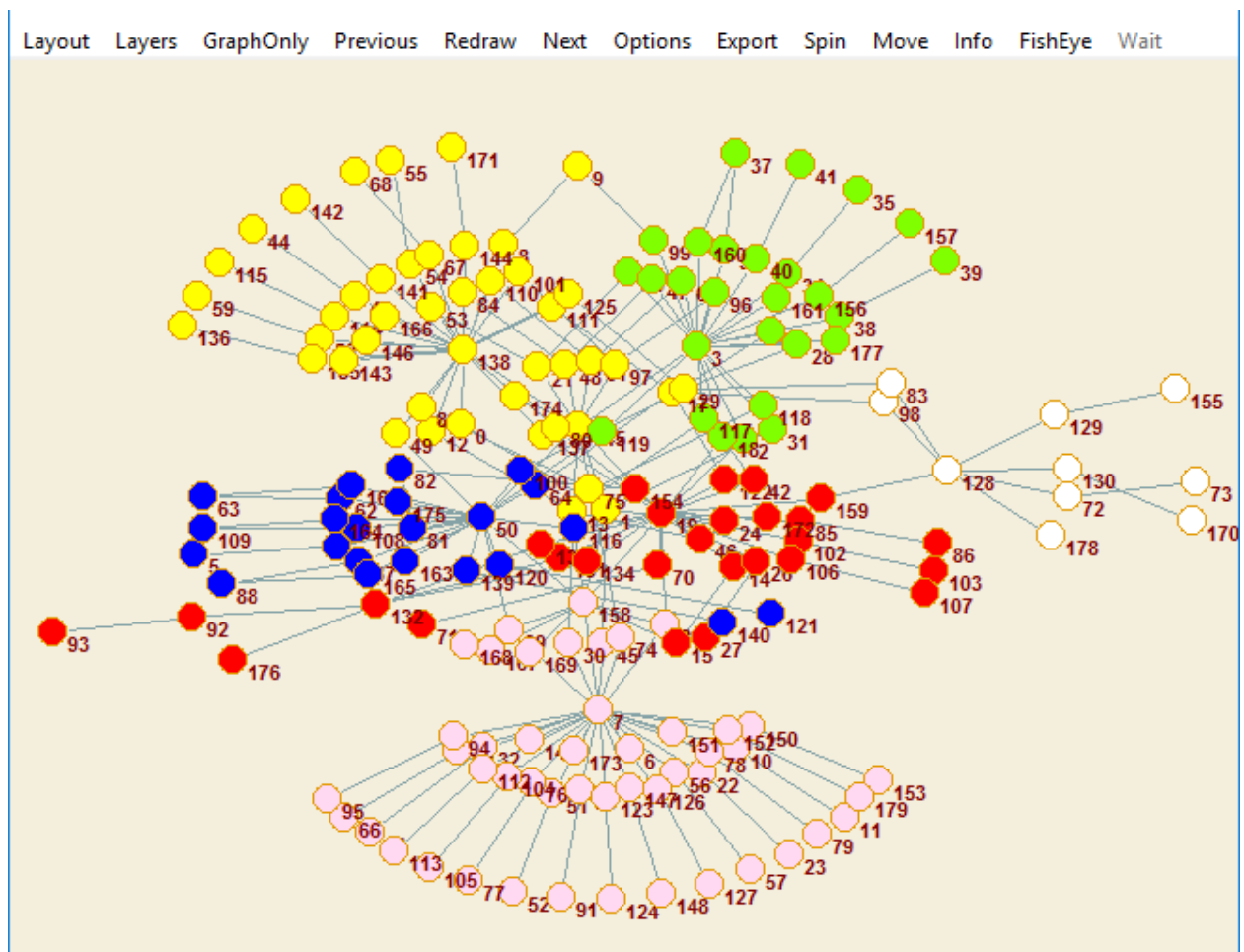
## Visualization Results for Datasets:

1. **AT&T Web Network**
   Total nodes – 180.
   Total edges – 228.
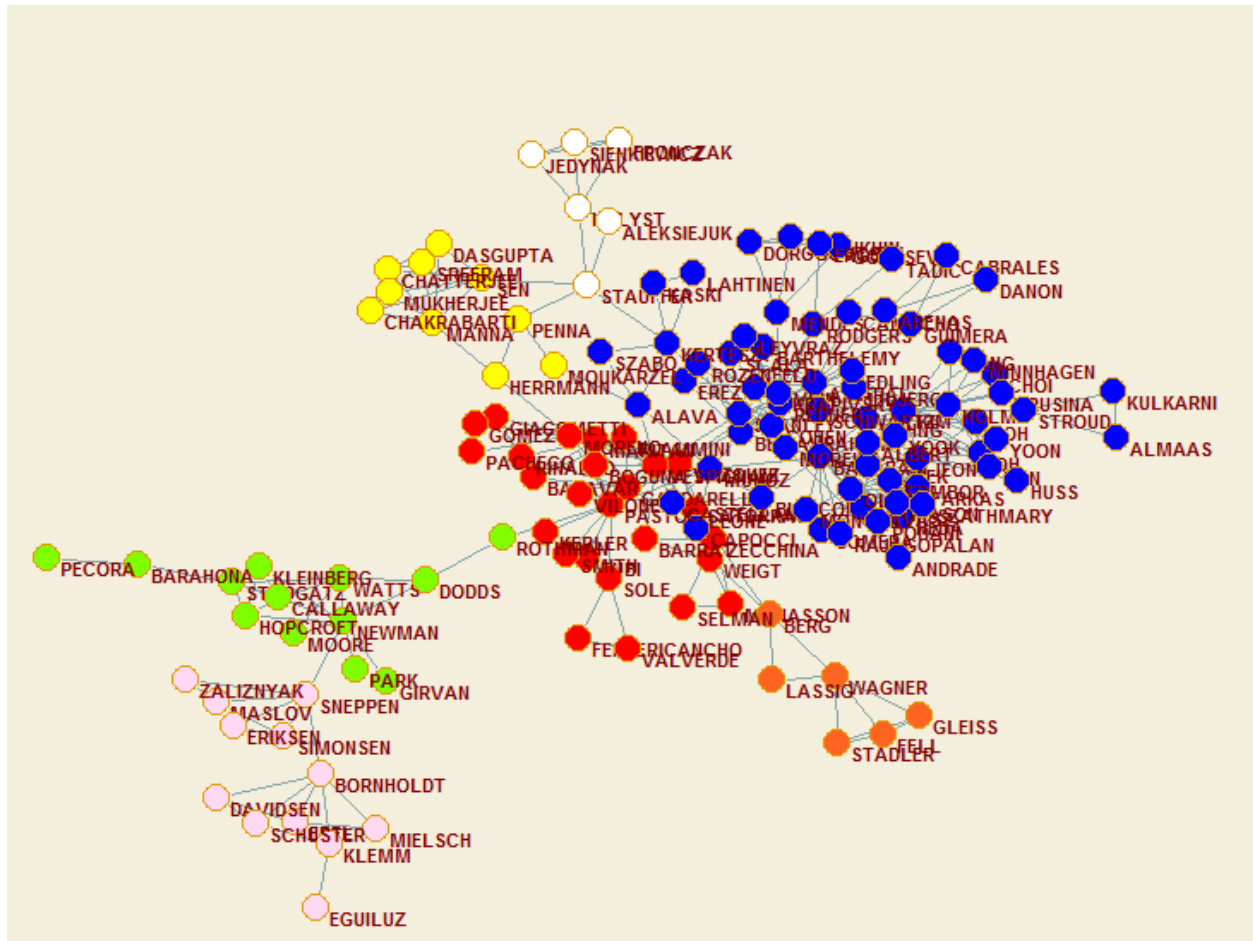   Total clusters - 6
   **Power – 30     Inflation - 30**

## 2. Physics Collaboration Network

Total number of nodes – 142

Total Edges – 340

Total Clusters – 7.

**Power: 30    Inflation: 50**

## 3. **Yeast Undirected Metabolic**

Total number of nodes – 359
Total number of edges – 435
Total number of cluster – 9

**Power - 140     Inflation - 140**