

# Assignment 5

CS532-s16: Web Sciences

Spring 2016

John Berlin

Generated on March 3, 2016

# 1

## Question

We know the result of the Karate Club (Zachary, 1977) split. Prove or disprove that the result of split could have been predicted by the weighted graph of social interactions. How well does the mathematical model represent reality?

Generously document your answer with all supporting equations, code, graphs, arguments, etc.

Useful sources include:

- \* Original paper

<http://aris.ss.uci.edu/~lin/76.pdf>

- \* Slides

<http://www-personal.umich.edu/~ladamic/courses/networks/si614w06/ppt/lecture18.ppt>

<http://clair.si.umich.edu/si767/papers/Week03/Community/CommunityDetection.pptx>

- \* Code and data

[http://networkx.github.io/documentation/latest/examples/graph/karate\\_club.html](http://networkx.github.io/documentation/latest/examples/graph/karate_club.html)

<http://nbviewer.ipython.org/url/courses.cit.cornell.edu/info6010/resources/11notes.ipynb>

<http://stackoverflow.com/questions/9471906/what-are-the-differences-between-community-detect>

<http://stackoverflow.com/questions/5822265/are-there-implementations-of-algorithms-for-commu>

<http://konect.uni-koblenz.de/networks/ucidata-zachary>

<http://vlado.fmf.uni-lj.si/pub/networks/data/ucinet/ucidata.htm#zachary>

<https://snap.stanford.edu/snappy/doc/reference/CommunityGirvanNewman.html>

[http://igraph.org/python/doc/igraph-pysrc.html#Graph.community\\_edge\\_betweenness](http://igraph.org/python/doc/igraph-pysrc.html#Graph.community_edge_betweenness)

## Answer

Before proving whether or not the result of split could have been predicted, let us first understand the problem. Zachary's Karate Club [4] is a very well known graph problem, where the social interactions between two groups "*Mr. Hi* and *John A*", caused a split in that social community. The club as seen in figure 1 is the original club with the results of the split shaded according to which faction each member went too. As seen in figure 1, the graph(club) is highly connected, which leaves us to wonder how would we figure out how to predict the division solely using the graph itself.

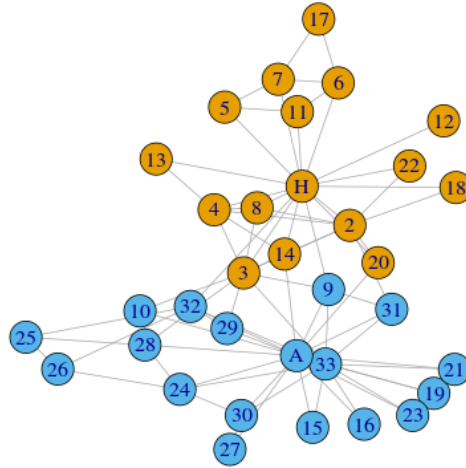


Figure 1: Zachary's Karate Club

To determine if the graph structure has enough information to predict the split, I used the edge betweenness algorithm as described by Girvan, Newman in their paper Community structure in social and biological networks [1]. In that paper they stated, that the "betweenness of an edge as the number of shortest paths between pairs of vertices that run along it". Furthermore they state that "edges connecting communities will have high edge betweenness" and "by removing those edges" the communities are exposed. Now to take from Dr. Nelson's lecture #6 "Edge betweenness: Total amount of flow an edge carries

between all pairs of nodes where a single unit of flow between two nodes divides itself evenly among all shortest paths between nodes". I believe that is straight from the Girvan Newman paper when they generalize Freeman's betweenness centrality[1].

This algorithm as stated in the paper by Girvan, Newman is

1. Calculate the betweenness for all edges in the network
2. Remove the edges with the highest betweenness
3. Recalculate betweenness for all edges affected by the removal
4. Repeat from step 2 until no edges remain

In using this algorithm to find the split, I must add one caveat to step number 4, which is to stop when the number of communities desired has been reached. To know when that stopping point has been reached is easy, simply count the number of connected components in the graph. The connected components at the start of the algorithm will be one, as the original graph is one piece. When the algorithm has removed enough edges the number connected components will be increased; in short I am counting the number of connected subgraphs [3]. I implemented this algorithm in Python using the Networkx library and can be seen in listing 1 produced the split seen in figure 3. Please see the comments for further details. But stated simply the program runs as such

1. Calculate the betweenness for all edges in the network
2. Get the maximum value from the dictionary gotten from `edge_betweenness_centrality`
3. Iterate through the returned dictionary and remove the edges whoes score is equal to max. This is to account for ties if any
4. Repeat until we reach the desired number of communities denoted `splitTo` by checking the number of connected components in the graph

The results of my implementation can be seen in table 1. The actor column represents the member number from the club with actor 1 being Mr. Hi and actor 34 being John A. The following two columns show the faction and result of the split from the Zachary paper [4]. The final two columns show the faction membership as perdicted my Girvan-Newman implementation and if it matched membership correctly, hit for yes and miss for no.

From the Zachary paper his model had a 97% percent accuracy rating missing only 1 membership out of the 33. Whereas my model missed 4 out of the 34 resulting in a 88% accuracy score. In comparison to the work done by Zachary I as well missed the faction that actor number 9 would go too. But it also must be noted that Girvan-Newman used the unweighted version of the Zachary graph.

This model as produced by the python implementation seemed to be off, so

I decided to use R to see if I did anything wrong. Why because as I went looking into graph libraries the igraph R package kept coming up so I decided to use it.

The code used in my R implementation can be found in listing 3 and the resulting graph split in figure 2. Now this implementation is the most interesting. I used the igraph library and consulted the pdf reference that accompanied it on cran[2]. In the library is this method called `cluster_edge_betweenness` which returns to us a communities object. This communities class performs the Girvan Newman algorithm for us and returns to us the edges removed, their betweenness score and much more. The comments in the code contain details on what is happening.

Since I already did the bigger analysis for my python implementation I will briefly break down the improvements. The r version correctly said the actor 9 would go to the Officers faction, actor number 14 was incorrectly placed in the Officers faction but actor 22 was correctly placed in Mr. Hi's faction. With these improvements I now miss only 2 out of the 34 resulting in a 94% accuracy score.

By this I can safely say that the split can be predicted with 90% accuracy of the actual split.

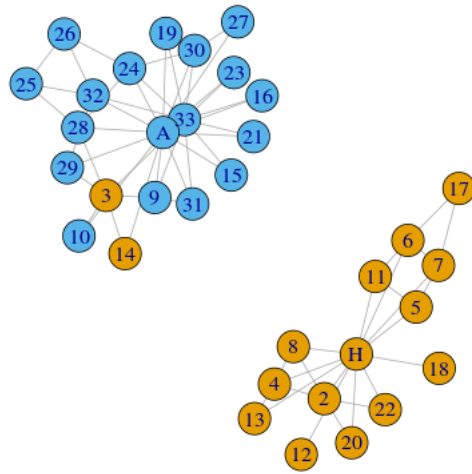


Figure 2: Club split by R implementation of Girvan Newman

```

1 import networkx as nx
2
3 def gnewman(club, splitTo = 2):
4     itteration = 0
5     # ok so why do I check the number of connected components
6     # for an undirected graph it is know that a connected component
7     # of an
8     # an undirected graph is a subgraph in which any two vertices
9     # are connected to each other by paths
10    # this is useful for this application since we are splitting a
11    # graph into two subgraphs
12    # ie to mathematically represent the splitting of the club
13    while nx.number_connected_components(club) < splitTo:
14        # returns to us edges with the weights
15        between = nx.edge_betweenness centrality(club, normalized=
16        False)
17        # we want the edges with the highest edge betweenness
18        centrality
19        # there might be ties so just get the max betweenness
20        m = max(between.values())
21        # unpack the tuple returned to us by between.items ((u,v),
22        maxBetweenScore)
23        for (hU,hV), val in between.items():
24            # check to see if m(max betweenness score) is equal to
25            val
26            # removes ties along the way
27            if val == m:
28                club.remove_edge(hU,hV)
29                print("removed edge %s—%s with betweenness score
30                of %f"%(hU,hV,m))
31                itteration += 1
32
33            print("—————")
34            # this print out can be uncommented it simply shows the
35            # same metric as described two different ways
36            # print(nx.number_connected_components(club),len(list(nx.
37            # connected_component_subgraphs(club))))
38            print("total iterations %d for splitting into %d"%(itteration,
39            splitTo))
40
41 if __name__ == "__main__":
42
43     for splitTo in range(2,6):
44         #this is a weighted version of the graph from http://www.
45         nexus.igraph.org/api/dataset_info?id=1&format=html
46         karateClub = nx.read_graphml("karate.GraphML")
47         # karateClub = nx.read_gml("nxKclub2.gml")
48         # karateClub = nx.karate_club_graph()
49         gnewman(karateClub, splitTo)
50         print("-----\n")
51         nx.write_dot(karateClub, "kc%d.dot"%splitTo)

```

Listing 1: Girvan Newman edge betweenness centrality

Actor	Faction Membership From Data	Zachary Membership From Split	Girvan- Newman Faction Membership From Split	Hit Miss for my Girvan- Newman
1	Mr. Hi	Mr. Hi	Mr. Hi	Hit
2	Mr. Hi	Mr. Hi	Mr. Hi	Hit
3	Mr. Hi	Mr. Hi	Officer	Miss
4	Mr. Hi	Mr. Hi	Mr. Hi	Hit
5	Mr. Hi	Mr. Hi	Mr. Hi	Hit
6	Mr. Hi	Mr. Hi	Mr. Hi	Hit
7	Mr. Hi	Mr. Hi	Mr. Hi	Hit
8	Mr. Hi	Mr. Hi	Mr. Hi	Hit
9	Mr. Hi	Officers	Officers	Miss
10	Officers	Officers	Mr. Hi	Miss
11	Mr. Hi	Mr. Hi	Mr. Hi	Hit
12	Mr. Hi	Mr. Hi	Mr. Hi	Hit
13	Mr. Hi	Mr. Hi	Mr. Hi	Hit
14	Mr. Hi	Mr. Hi	Mr. Hi	Hit
15	Officers	Officers	Officers	Hit
16	Officers	Officers	Officers	Hit
17	Mr. Hi	Mr. Hi	Mr. Hi	Hit
18	Mr. Hi	Mr. Hi	Mr. Hi	Hit
19	Officers	Officers	Officers	Hit
20	Mr. Hi	Mr. Hi	Mr. Hi	Hit
21	Officers	Officers	Officers	Hit
22	Mr. Hi	Mr. Hi	Officers	Miss
23	Officers	Officers	Officers	Hit
24	Officers	Officers	Officers	Hit
25	Officers	Officers	Officers	Hit
26	Officers	Officers	Officers	Hit
27	Officers	Officers	Officers	Hit
28	Officers	Officers	Officers	Hit
29	Officers	Officers	Officers	Hit
30	Officers	Officers	Officers	Hit
31	Officers	Officers	Officers	Hit
32	Officers	Officers	Officers	Hit
33	Officers	Officers	Officers	Hit
34	Officers	Officers	Officers	Hit

Table 1: Comparisons of the results of split as generated by my Girvan Newman implementation to Zachary’s predictions and the actual data as was done in the original paper



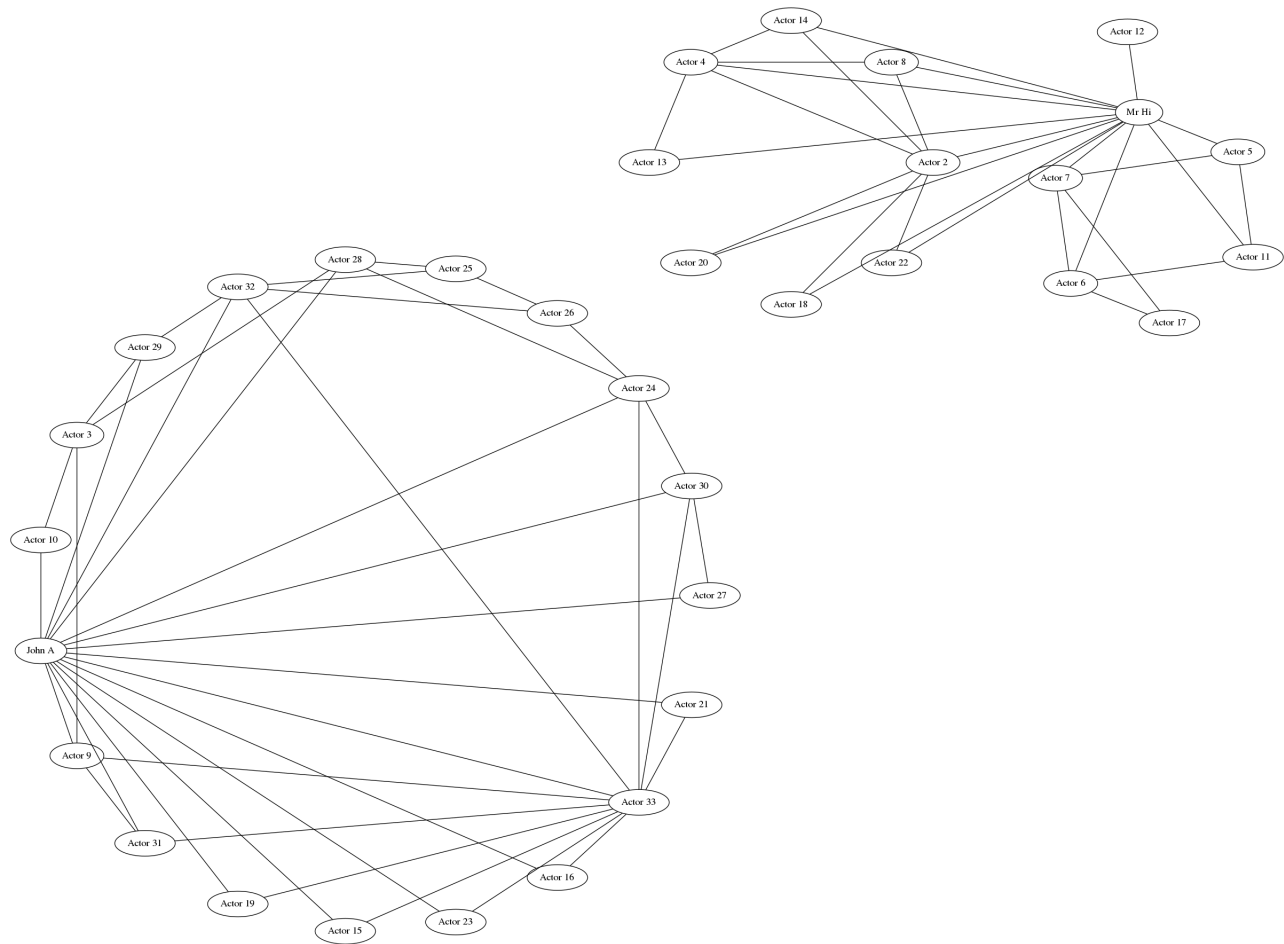


Figure 3: Club split by my implementation of Girvan Newman

```
1 $ circo -Tpng -o kc.png kc2.dot  
2 $ circo -Tpng -o kc3.png kc3.dot  
3 $ circo -Tpng -o kc4.png kc4.dot  
4 $ circo -Tpng -o kc5.png kc5.dot
```

Listing 2: Convert dot files to png

```

1 library(igraph)
2 library(igraphdata)
3
4 #ok so lets do this in R
5 data(karate)
6 karateClub <- karate
7 #mmm this looks interesting
8 #edge.betweenness.community(cluster_edge_betweenness),65 from the
   pdf on cran
9 ceb <- cluster_edge_betweenness(karateClub)
10 #what does this give us a communities class
11 print(class(ceb))
12 #look what we have here! it removed edges for us well sorta
13 print(ceb$removed.edges)
14 #this will print edge betweenness ;
15 print(ceb$algorithm)
16
17 #set to true to see debug output
18 debug <- TRUE
19
20 #how many communities do we wish to split to
21 splitTo <- 2
22
23 #ok some set up
24 #want to preserve the final graph that we create
25 # I will not be touching the original karate club
26 finalPlot <- karateClub
27 plot.igraph(finalPlot)
28
29 #ok so I want have a counter to keep track
30 #of where I am in the algorithm
31 removeIndex <- 1
32 #use the same scheme as python I want to make sure
33 #the number of connected componets are less than 2
34 #for the Zachary Analysis
35 numConnected <- count_components(karateClub)
36
37 #loop
38 while(numConnected < splitTo){
39   #ok so this feels like cheating here
40   #the scheme is to progressively remove edges the ebc communities
     class calculated for us
41   #if you start at index 0 you get nothing so I start at 1 to
     removeIndex -1
42   #base removeIndex=1 seq(1,removeIndex-1) yeild 1 0
43   #so the ebc$removed.edges[seq(1,removeIndex-1) ] yeild edge at
     pos 1 :)
44   deletedG <- delete.edges(karateClub, ceb$removed.edges[seq(1,
     removeIndex-1)])
45   #deletedG is now the result of iteration N of Girvan Newman
46
47   #increment our iteration of Girvan Newman up by one
48   removeIndex <- removeIndex + 1
49   # plot.igraph(deletedG)
50   if(debug){
51     print(seq(1,removeIndex-1))
52     print(ceb$edge.betweenness[seq(1,removeIndex-1)])

```

```

53 |     print(ceb$removed.edges[seq(1,removeIndex-1)])
54 |     print("+++++")
55 | }
56 | #get the number of connected componets ie
57 | #have we split the graph into two subgraphs
58 | numConnected <- count_components(deletedG)
59 | finalPlot <- deletedG
60 | }
61 |
62 | plot.igraph(finalPlot)

```

Listing 3: Girvan Newman edge betweenness centrality in R

## 2

### Question

2. We know the group split in two different groups. Suppose the disagreements in the group were more nuanced -- what would the clubs look like if they split into groups of 3, 4, and 5?

### Answer

Using the same R code as in the second implementation I changed the splitTo value to 3 communities figure 4, 4 communities figure5, and 5 communities figure6. Those figures are seen below.

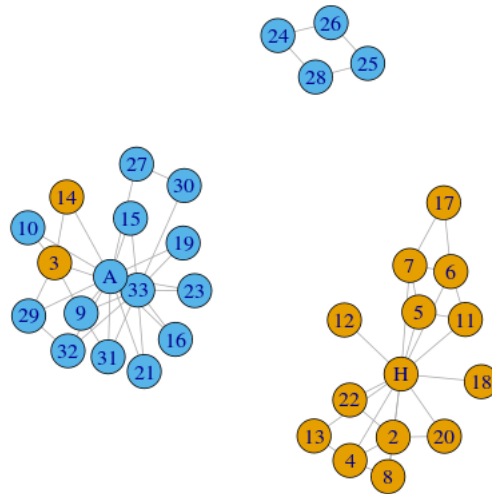


Figure 4: Karate Club Graph Split Into 3 Communities

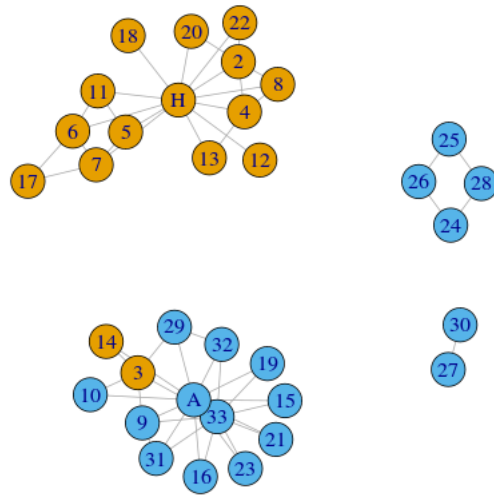


Figure 5: Karate Club Graph Split Into 4 Communities

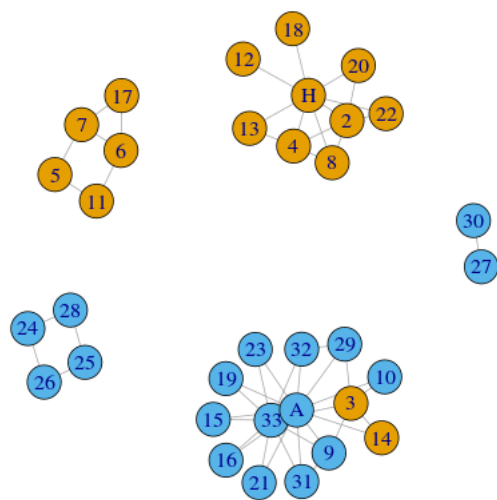


Figure 6: Karate Club Graph Split Into 5 Communities

## References

- [1] GIRVAN, M., AND NEWMAN, M. E. J. Community structure in social and biological networks. *arXiv:cond-mat/0112110* (2001).
- [2] MANY SEE AUTHORS FILE. igraph: Network analysis and visualization, 2015.
- [3] WIKIPEDIA. Connected component (graph theory) — wikipedia, the free encyclopedia, 2015. [Online; accessed 3-March-2016].
- [4] ZACHARY, W. W. An information flow model for conflict and fission in small groups. *Journal of anthropological research* (1977).