

1. $Q = \sum_r (e_{rr} - a_r^2)$ Eq. (7.76)

e_{rr} is given by the diagonals of the table given in Problem 1, since this value represents the fraction of all edges that connect like vertex to like vertex (in other words, a given edge e that is included in a diagonal has one endpoint connected to a man of ethnicity r , and the other to a woman of ethnicity r). a_r , the number of edges that connect to a given vertex type r , can be found by: $A_{rr} + \sum_{s \neq r} (A_{rs} + A_{sr})$. This represents the number of edges with at least one endpoint connecting to vertex type r .

Ethnicity	e_{rr}	a_r	a_r^2
Black	0.258	0.352	0.124
Hispanic	0.157	0.292	0.085
White	0.306	0.494	0.244
Other	0.016	0.119	0.014

This gives us the following values:

So, by Eq.(7.76), $Q = 0.27$

Overall, the community exhibits distinct homophily. The only outlier is the *other* vertex type, which is not surprising, considering that this vertex type is actually a collection of vertices of varying types. These smaller subtypes may be small minorities of the community as a whole.

2. a. A line graph with n vertices ($n-1$ total edges, $2n-2$ total endpoints), whose nodes are divided into two contiguous groups, r and $\neg r$ results in the following configuration:

Component	(e_{rr})	(a_r)
r	$\frac{r-1}{n-1}$	$\frac{2r-1}{2n-2}$
$\neg r$	$\frac{n-r-1}{n-1}$	$\frac{2n-2r-1}{2n-2}$

e_{rr} is the number of nodes that link vertices of the same component.

a_r is calculated by $\frac{\text{number of endpoints attached to a component's vertices}}{\text{total number of endpoints in the graph}}$

$$Q = \sum_r (e_{rr} - a_r^2)$$

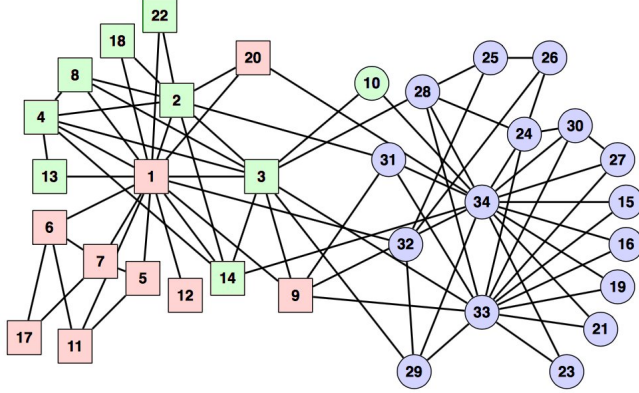
$$Q = \frac{r-1}{n-1} - \left(\frac{2r-1}{2n-2}\right)^2 + \frac{n-r-1}{n-1} - \left(\frac{2n-2r-1}{2n-2}\right)^2$$

$$Q = \frac{3-4n+4rn-4r^2}{2(n-1)^2}$$

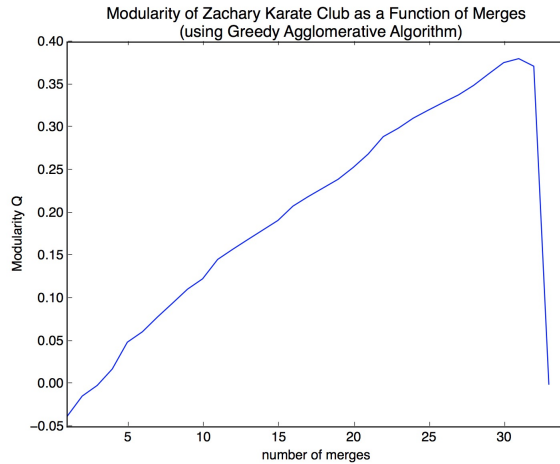
2. b. If the graph has n nodes where n is even, then each component has $r = \neg r = \frac{n}{2}$ nodes. To maximize Q , we can use the first and second derivatives of the result from 2.a. $\frac{d}{dr} \frac{3-4n+4rn-4r^2}{2(n-1)^2} = \frac{2(n-2r)}{(n-1)^2}$ which equals 0 at $r = \frac{n}{2}$, and $\frac{d^2}{du^2} = \frac{-4}{(n-1)^2}$, which is negative for

all values of r , so $r = \frac{n}{2}$ is a global maximum for the equation $Q = \frac{3-4n+4rn-4r^2}{2(n-1)^2}$. This tells us that, when n is even, the division resulting in optimal modularity is one where both components have an equal number of nodes.

3. Using the greedy agglomerative algorithm on the Zachary Karate Club network, the maximum Q score achieved was 0.3806 with three groups (after 31 merges), as shown here:



Q as a function of number of merges gives:



The Normalized Mutual Information is $\frac{2I(C,C')}{H(C)+H(C')}$

where $H(C) = -\frac{17}{34}\log_2\frac{17}{34} - \frac{17}{34}\log_2\frac{17}{34} = 1$

$H(C') = -\frac{8}{34}\log_2\frac{8}{34} - \frac{9}{34}\log_2\frac{9}{34} - \frac{17}{34}\log_2\frac{17}{34} = 1.498751$

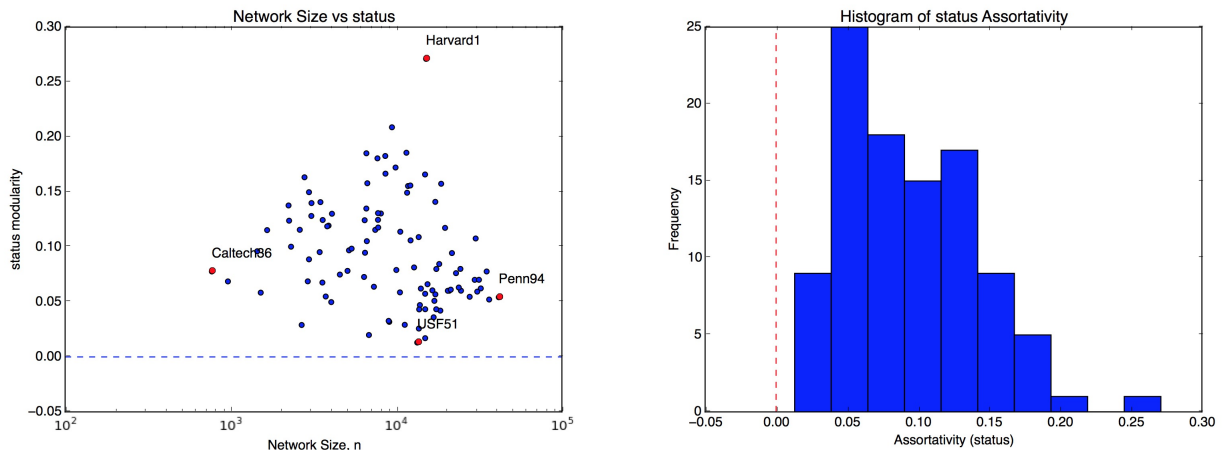
and $I(C, C') = H(C) - H(C|C') = 1 - (\frac{8}{34}\log_2 1 - \frac{8}{34}\log_2\frac{8}{9} - \frac{1}{34}\log_2\frac{1}{9} - \frac{1}{34}\log_2\frac{1}{17} - \frac{16}{34}\log_2\frac{16}{17})$

$= 1.0 - .0.294594 = 0.705406$

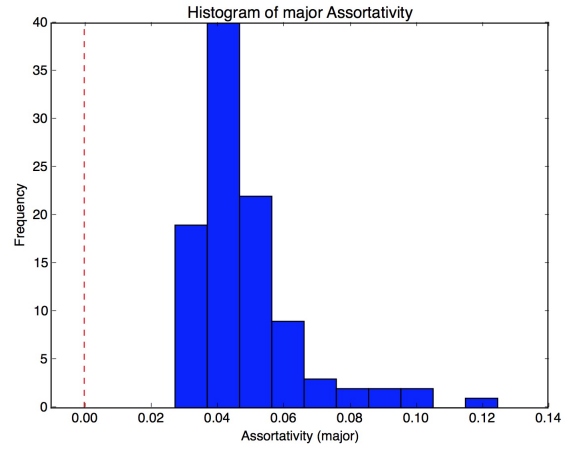
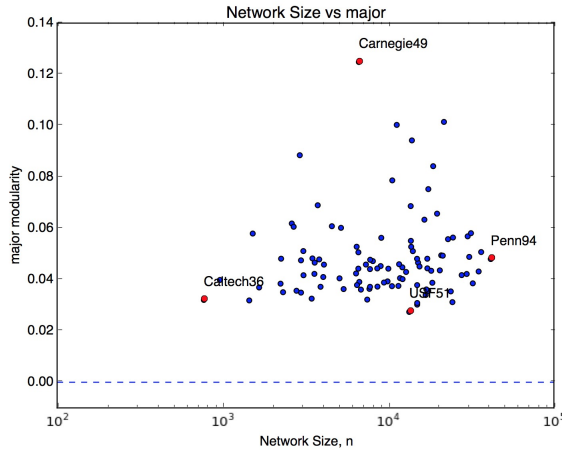
$NMI = \frac{2(0.705406)}{2.498751} = 0.564607$

Since $NMI = 0.564607$, this would suggest that, without knowing the 'truth on the ground', it is possible to make relatively accurate, but not perfect, predictions about group membership with the greedy agglomerative algorithm. It is interesting that the greedy algorithm gave three partitions, while the truth on the ground was made up of only two. Nodes 2 and 3 seem to be the driving cause of this added cluster, forming many triangles with each other and surrounding nodes, with node 1 doing the same for the remaining nodes in the original (red) partition.

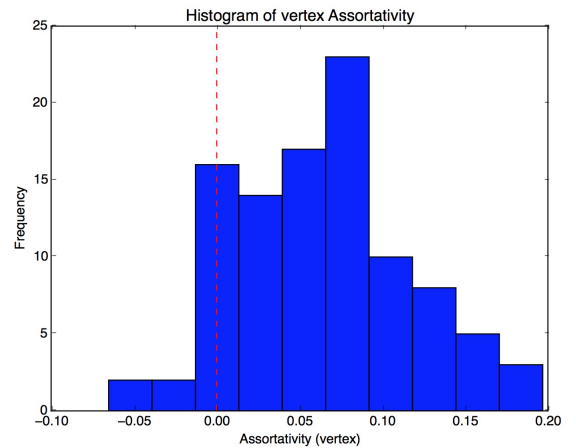
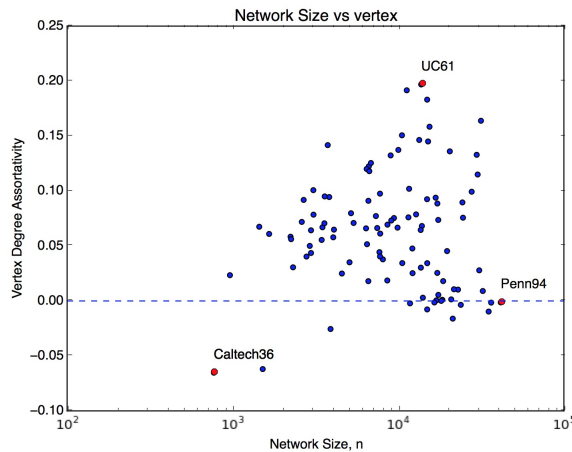
4. Assortativity/modularity was calculated using the iGraph library for Python, as the modularity and assortativity functions therein are calculated using the same methods as presented in class and lecture notes.



Network size and modularity of status show distinct homophily. Networks of all sizes exhibit positive modularity with respect to status (faculty, student, alumni, etc). The most common degree of assortativity is around the 0.05 range, but the histogram has a rather long, positive tail. This seems logical, since current students most likely connect with other students, faculty to faculty, etc. Many universities most likely have a policy regarding faculty/student interactions outside of official business, so this may act to reinforce the homophily of the networks by actively discouraging connections between some groups.



Networks also exhibit homophily with regards to major. The distribution is considerably more narrow than seen in the status data, with 40% of networks having assortativity with respect to major at around 0.045. Again, this is not surprising, considering that students most likely spend a lot of time with other students of the same major. It is impossible to tell from the data given if these are purely 'friendship' connections, or if the data reflects a more academic connection between nodes, such as project or study groups.

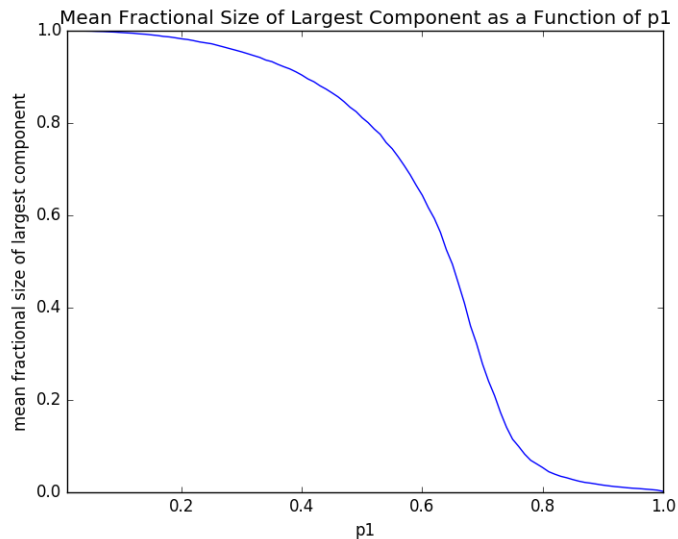


Assortativity with respect to vertex degree exhibits what seems to be the widest range of assortativity of all data sets examined, with some universities showing negative assortativity. However, the histogram shows a wide distribution from 0.0 to 0.75, and with a long positive tail. The lower vertex degree associations could be caused by the friendship paradox, as discussed earlier in class. Since the average person's friends have more friends than she does, it follows that nodes link to nodes with a dissimilar degree number. This does not explain the wide distribution, however.

6. To compute fractional sized of configuration model graphs, the networkx library for Python was used for convenience, as it has built-in functions for the generation of configuration model random graphs and finding the connected components therein.

Part 1: To determine the mean fractional size of the largest component for a network with $n = 10^4$ vertices, and with $p_1 = 0.6$ and $p_3 = 1 - p_1$, 2663 iterations were performed. Each iteration consisted of the generation of a random graph by the `networkx.configuration_model()` function, using randomly generated degree sequences, matching the probabilities given. The number 2663 was chosen to achieve a 99% confidence level and a confidence interval of 2.5^* . Under these conditions, the mean fractional size of the specified random network was 0.64204.

Part 2: Again using the networkx library, 666 iterations were performed for each p_1 value, from $p_1 = 0.01$ to $p_1 = 1.0$. 666 iterations were done to achieve a confidence level of 99% and a confidence interval of 5^* , giving the following results:



The phase change, where the giant component disappears, happens as p_1 grows larger than 0.6 .

* per <http://www.surveysystem.com/sscalc.htm>

```

1  # Author: Donovan Guelde
2  # CSCI 5352 PS3 Question 3
3  # references: online documentation for numpy,
4  #             http://stackoverflow.com/questions/9754729/remove-object-from-a-list-of-objects-in-python
5  # Collaborators: None
6
7  import numpy as np
8  import os
9  import matplotlib.pyplot as plt
10 import networkx as nx
11
12 #fileName="test.txt"
13 fileName="../../CSCI5352_Data/karate_club_edges.txt"
14
15 def inModGroup(i,maxModularityGroups):
16     if i in maxModularityGroups[0]:
17         return 0
18     if i in maxModularityGroups[1]:
19         return 1
20     if i in maxModularityGroups[2]:
21         return 2
22
23 class Network:
24     def __init__(self,fileName):
25         self.groups = [] #empty array at instantiation, filled and updated as merges are performed
26         self.n = 0 #number of nodes in network
27         self.associationMatrix = self.readFile(fileName) #simple graph
28         self.m = np.sum(self.associationMatrix)/2 #number edges
29         self.regularization = (2*float(self.m)) #so we only have to calculate it once...
30         self.eMatrix = self.get_e_matrix()
31     def readFile(self,fileName):
32         with open(fileName,'r') as f: #make 2d numpy array of appropriate size
33             temp=-1
34             lastNode=-1
35             for line in f:
36                 line = line.split()
37                 temp = max(int(line[0]),int(line[1]))
38                 if (temp>lastNode):
39                     lastNode=temp #finds the highest numbered node
40             associationMatrix = np.zeros((lastNode,lastNode))
41             self.n = int(lastNode) #assumes no gaps in node labelling
42             f.seek(0,0)
43             lines = f.readlines()
44             for line in lines:
45                 line = line.split()
46                 associationMatrix[int(line[0])-1][int(line[1])-1] = 1 #make it undirected...
47                 associationMatrix[int(line[1])-1][int(line[0])-1] = 1
48             for index in range(0,self.n): #self.groups is a list of lists
49                 self.groups.append([]) #add empty list for every vertex in graph
50                 self.groups[index][0] = index+1 #place every vertex in its own group
51         f.close()
52         return (associationMatrix)
53
54     def inGroup(self,i): #returns group that node i belongs to

```

```

55         group=0
56         node = i
57         for index in range (0,len(self.groups)):
58             if (node in self.groups[index]):
59                 group=index
60                 break
61         return group
62
63     def get_e_matrix(self): #updates e matrix (used after merge is performed)
64         numberGroups=0
65         for index in range(0,len(self.groups)):
66             if (self.groups[index][0]): #if group has member/members
67                 numberGroups+=1
68         eMatrix = np.zeros((numberGroups,numberGroups))
69         for r in range (0,numberGroups):
70             for s in range (0,numberGroups):
71                 tempSum=0.0
72                 for i in range (0,self.n):
73                     if(self.inGroup(i+1) == r):
74                         for j in range(0,self.n):
75                             if (self.associationMatrix[i][j]==1):
76                                 if (self.inGroup(j+1) == s):
77                                     tempSum+=1
78                 if (tempSum!=0):
79                     eMatrix[r][s]=(tempSum/self.regularization)
80         return eMatrix
81
82     def findDeltaQ(self,u,v): #returns delta Q between groups u and v
83         a_u = np.sum(self.eMatrix[u])
84         a_v = np.sum(self.eMatrix[v])
85         return (2*(self.eMatrix[u][v]-(a_u*a_v)))
86
87     def findGreatestDeltaQ(self): #returns (greatest delta Q,
88         # index of first group to merge, index of second group to merge)
89         deltaQ = (float('-inf'),0,0)
90         for index in range(0,len(self.groups)):
91             for index2 in range (index+1,len(self.groups)):
92                 temp = self.findDeltaQ(index,index2)
93                 if (deltaQ[0] < temp):
94                     deltaQ = (temp,index,index2)
95         return deltaQ
96
97     def mergeGroups(self,r,s): #merge groups r and s into r, delete s
98         r=int(r)
99         s=int(s)
100         for item in self.groups[s]:
101             self.groups[r].append(item)
102         del self.groups[s]
103         return
104
105     def getQ(self): #returns Q of network
106         sum=0
107         for index in range(0,len(self.groups)):
108             sum2=0
109             for index2 in range(0,len(self.groups)):
110                 sum2+=self.eMatrix[index][index2]
111             sum+=self.eMatrix[index][index]-pow(sum2,2)
112         return sum

```

```

113
114 def main():
115     graph = Network(fileName)
116     counter=0
117     maxModularity=float('-inf')
118     modularity=np.zeros((graph.n))
119     while (len(graph.groups)>1):
120         temp=graph.findGreatestDeltaQ() #temp is a triple, (max delta Q, group to merge, group to merge)
121         Q = graph.getQ()
122         graph.mergeGroups(temp[1],temp[2])
123         modularity[counter] = Q
124         print Q
125         print graph.groups
126         counter+=1
127         graph.eMatrix=graph.get_e_matrix()
128         if (maxModularity<graph.getQ()):
129             maxModularity=graph.getQ()
130             length = len(graph.groups)
131             maxModularityGroups=[[]]*length
132             #print graph.getQ()
133             for index in range(0,length):
134                 maxModularityGroups[index] = graph.groups[index]
135     print maxModularity
136     print maxModularityGroups
137
138     plt.ylabel('Modularity Q')
139     plt.xlabel('number of merges')
140     plt.xlim(1,34)
141     plt.title('Modularity of Zachary Karate Club as a Function of Merges\n (using Greedy Agglomerative Algorithm)')
142     plt.plot(modularity)
143     plt.savefig("karateMerge.jpg")
144
145     main()

1  # Author: Donovan Guelde
2  # CSCI 5352 PS3 Question 4
3  # references: online documentation for numpy and igraph
4  # Collaborators: None
5
6  import numpy as np
7  import os
8  import igraph
9
10 class Network:
11     def __init__(self,fileName):
12         self.n = 0
13         self.gender=[]
14         self.status=[]
15         self.major = []
16         self.vertexDegree=[]
17         self.associationMatrix = self.readFile(fileName)
18         self.g = igraph.Graph.Adjacency((self.associationMatrix == 1).tolist())
19         self.m = np.sum(self.associationMatrix)/2.0
20         self.regularization = 1/(2*self.m) #caluclate once
21     def readFile(self,fileName):
22         with open("./facebook100txt/"+fileName+"_attr.txt",'r') as f: #get n and attributes from _attr.txt
23             counter=0
24             for line in f:

```



```

25         counter+=1
26         self.n = counter-1
27         associationMatrix = np.zeros((self.n,self.n))
28         self.gender = [0]*self.n #arrays to track gender, status, major and degree of vertexes
29         self.status = [0]*self.n
30         self.major = [0]*self.n
31         self.vertexDegree = [0]*self.n
32         f.seek(0,0)
33         f.next() #skip the label row
34         counter=0
35         for line in f: #populate the attribute arrays
36             line = map(int,line.split())
37             self.gender[counter] = int(line[2]) #gender of vertexes where index=vertex
38             self.status[counter] = int(line[1]) #ditto...
39             self.major[counter] = int(line[3])
40             counter+=1
41     f.close()
42
43     with open("./facebook100txt/"+fileName+".txt",'r') as f: #construct association matrix
44         lines = f.readlines()
45         for line in lines:
46             line = line.split()
47             associationMatrix[int(line[0])-1][int(line[1])-1] = 1
48             associationMatrix[int(line[1])-1][int(line[0])-1] = 1 #make it undirected
49     f.close()
50     for index in range(0,self.n): #populate the vertex degree array,
51         self.vertexDegree[index] = np.sum(associationMatrix[index],axis=0)
52     return associationMatrix
53
54     def getQ(self,attribute): #returns Q of network
55         if attribute == "gender":
56             membership = self.gender
57         if attribute == "status":
58             membership = self.status
59         if attribute == "major":
60             membership = self.major
61         Q = igraph.Graph.modularity(self.g,membership)
62         return Q
63
64     def calculateAssortativity(self):
65         assortativityCoefficient=igraph.Graph.assortativity(self.g,self.vertexDegree)
66         return assortativityCoefficient
67
68     def main():
69         plotArray=np.empty((100,2)) #an array of points to plot
70         nameArray=[""]*100 #array to hold names of schools where [index] corresponds to plotArray[index]
71         nextUniversity = [2]
72         lastUniversity = [2]
73         genderModularity=[""]*100
74         statusModularity=[""]*100
75         majorModularity=[""]*100
76         vertexAssortativity=[""]*100
77         names=[""]*100
78         nValues = [""]*100
79         counter=0
80         for file in os.listdir("./facebook100txt/"):
81             if (file != ".DS_Store"):
82                 nextFile, fileExtension = os.path.splitext(file)

```

```

83         nextUniversity = nextFile.split('_')
84         if (str(nextUniversity[0]) != str(lastUniversity[0])):
85             nextGraph = Network(nextUniversity[0])
86             names[counter]=nextUniversity[0]
87             nValues[counter]=nextGraph.n
88             genderModularity[counter]=nextGraph.getQ("gender")
89             statusModularity[counter]=nextGraph.getQ("status")
90             majorModularity[counter]=nextGraph.getQ("major")
91             vertexAssortativity[counter]=nextGraph.calculateAssortativity()
92             nameArray[counter] = str(nextUniversity[0])
93             lastUniversity=nextUniversity
94             counter+=1
95         with open("./results/genderModularity.txt","w") as f:
96             np.savetxt(f,genderModularity,fmt='%s')
97         f.close()
98         with open("./results/statusModularity.txt","w") as f:
99             np.savetxt(f,statusModularity,fmt='%s')
100        f.close()
101        with open("./results/majorModularity.txt","w") as f:
102            np.savetxt(f,majorModularity,fmt='%s')
103        f.close()
104        with open("./results/vertexAssortativity.txt","w") as f:
105            np.savetxt(f,vertexAssortativity,fmt='%s')
106        f.close()
107        with open("./results/names.txt","w") as f:
108            np.savetxt(f,names,fmt='%s')
109        f.close()
110        with open("./results/nValues.txt","w") as f:
111            np.savetxt(f,nValues,fmt='%s')
112        f.close()
113
114    main()

1  # Donovan Guelde
2  # CSCI 5352
3  # EC #6
4  # references: networkx online documentation, https://docs.python.org/2/library/functions.html#max
5
6  import numpy as np
7  import networkx as nx
8  import matplotlib.pyplot as plt
9  import random
10
11  DOPART1=1
12  DOPART2=1
13  #calculate mean fractional size per part 1, Q6
14  if (DOPART1):
15      NUMBERITERATIONS=2663 #confidence level-99, confidence interval-2.5
16      resultsArray = np.zeros(NUMBERITERATIONS)
17      for index in range(0,NUMBERITERATIONS):
18          degreeSequence=[0]*1000
19          degreeSequence[0]=1
20          while (not nx.is_valid_degree_sequence(degreeSequence)):
21              for index2 in range (0,1000): #generate degree sequence
22                  rollTheDice=random.random()
23                  if (rollTheDice<=.6):
24                      degreeSequence[index2]=1
25                  else:

```

```

26         degreeSequence[index2]=3
27         graph = nx.configuration_model(degreeSequence)
28         resultsArray[index] = len(max(nx.connected_component_subgraphs(graph),key=len))
29     print "mean of largest component after",NUMBERITERATIONS,"iterations =",np.mean(resultsArray)
30     print "mean fractional size =",np.mean(resultsArray)/1000
31
32
33
34     if (DOPART2):
35         NUMBERITERATIONS=666 #confidence level-99, confidence interval-5
36         resultsArray = np.zeros(100)
37         p1Values = np.zeros(100)
38         p1=.01
39         counter=0
40         while (p1<=1.01):
41             iterationResultsArray = np.zeros(NUMBERITERATIONS)
42             for index in range(0,NUMBERITERATIONS):
43                 degreeSequence=[0]*1000
44                 degreeSequence[0]=1
45                 while (not nx.is_valid_degree_sequence(degreeSequence)):
46                     for index2 in range (0,1000): #generate degree sequence
47                         rollTheDice=random.random()
48                         if (rollTheDice<=p1):
49                             degreeSequence[index2]=1
50                         else:
51                             degreeSequence[index2]=3
52                 graph = nx.configuration_model(degreeSequence)
53                 iterationResultsArray[index] = len(max(nx.connected_component_subgraphs(graph),key=len))
54             resultsArray[counter]=np.mean(iterationResultsArray)/1000
55             p1Values[counter]=p1
56             p1+=.01
57             counter+=1
58
59         plt.ylabel('mean fractional size of largest component')
60         plt.xlabel('p1')
61         plt.title('Mean Fractional Size of Largest Component as a Function of p1')
62         plt.xlim(.01,1)
63         plt.plot(p1Values,resultsArray)
64         plt.savefig("meanFractionalSize.png")
65         plt.close()
66
67
68
69
70

```

```

1  # Donovan Guelde
2  # CSCI 5352 PS3
3  # Plotter for PS3, Q4
4
5  import numpy as np
6  import matplotlib.pyplot as plt
7
8  def readData(fileName):
9      counter=0
10     with open(fileName,'r') as f:
11         for line in f:
12             counter+=1

```

```

13         f.seek(0,0)
14         data = np.zeros((counter)) # array to hold data
15         counter=0
16         maxDataPoint = float('-inf')
17         maxIndex=0
18         minDataPoint = float('inf')
19         minIndex=0
20         for line in f:
21             data[counter] = float(line)
22             if data[counter] > maxDataPoint:
23                 maxDataPoint=data[counter]
24                 maxIndex=counter
25             if data[counter] < minDataPoint:
26                 minDataPoint=data[counter]
27                 minIndex=counter
28             counter+=1
29         return data, int(maxIndex), int(minIndex)
30
31     def readNames(fileName):
32         counter=0
33         with open(fileName,'r') as f:
34             for line in f:
35                 counter+=1
36         f.seek(0,0)
37         data = [""]*counter
38         counter=0
39         for line in f:
40             data[counter] = line
41             counter+=1
42         return data
43     def getData(attribute):
44         if attribute == "major":
45             data, maxPoint, minPoint = readData("./results/majorModularity.txt")
46         if attribute == "status":
47             data, maxPoint, minPoint = readData("./results/statusModularity.txt")
48         if attribute == "vertex":
49             data, maxPoint, minPoint = readData("./results/vertexAssortativity.txt")
50         return data,int(maxPoint),int(minPoint)
51
52
53
54     def main():
55
56         nArray,nMax,nMin = readData("./results/nValues.txt")
57         names = readNames("./results/names.txt")
58         attributes = ["major","status","vertex"]
59
60         for item in attributes:
61             # scatter plots
62             data,maxPoint,minPoint = getData(item)
63             if item == "vertex":
64                 plt.ylabel('Vertex Degree Assortativity')
65             else:
66                 plt.ylabel(item+" modularity")
67             plt.xlabel('Network Size, n')
68             plt.title('Network Size vs '+item)
69             plt.scatter(nArray,data)
70             if item=='major':

```

```

71         plt.ylim(-.01,.14)
72     plt.text(float(nArray[maxPoint]),float(data[maxPoint]),names[maxPoint])
73     plt.plot(float(nArray[maxPoint]),float(data[maxPoint]),'o',mfc='red')
74     plt.text(float(nArray[minPoint]),float(data[minPoint]),names[minPoint])
75     plt.plot(float(nArray[minPoint]),float(data[minPoint]),'o',mfc='red')
76     plt.text(float(nArray[nMax]),float(data[nMax]),names[nMax])
77     plt.plot(float(nArray[nMax]),float(data[nMax]),'o',mfc='red')
78     plt.text(float(nArray[nMin]),float(data[nMin]),names[nMin])
79     plt.plot(float(nArray[nMin]),float(data[nMin]),'o',mfc='red')
80     plt.axhline(0,linestyle='dashed')
81     plt.xscale('log')
82     plt.savefig(item+".jpg")
83     plt.clf()
84     plt.close()
85
86     # histograms
87     plt.hist(data)
88     if item=='status':
89         plt.xlim(-.05,.30)
90     if item=='major':
91         plt.xlim(-.01,.14)
92     plt.title('Histogram of '+item+' Assortativity')
93     plt.xlabel('Assortativity ('+item+')')
94     plt.ylabel('Frequency')
95     plt.axvline(0,linestyle='dashed', color='red')
96     plt.savefig(item+"Hist.jpg")
97     plt.clf()
98     plt.close()
99
100 main()

```