

Working at Google - good idea?

You meet Sergey Brin

You've been hired at Google to run their **Life Sciences division**, congratulations! On your first day after your move to California's Silicon Valley, **Sergey Brin** comes in to meet the new wizard pythonista. He gets talking and admits to you there's so many googlers that he just cannot keep track of all of them. But he would like to build a rolodex of **central** googlers. Not just central in terms of titles or degrees or awards, but central in terms of social connections. Sergey wants you to find out who the most *social* googlers are.



Sergey has an idea: He gives you a list of projects that googlers have worked on, a list of Cruises googlers have gone on together, and a list of sport clubs users belong to.

You get 3 excel files:

- 1) `gprojects` are the projects googlers have worked on. The first column is the name of google employees, the subsequent columns are the projects that they have worked on. Googlers love codes, so the projects are encoded as sentences of **four** words each. The sentences don't need to make sense.
- 2) `gcruises` are the cruises that googlers have gone on. each row has the name of the google employee in the first column, and then the cruises she has been on. Cruises are encoded as sentences of **three** words each.
- 3) `gsports` are the sport clubs that googlers belong to. The first column contains googler employee name, subsequent columns the sport clubs they belong to. Sport clubs are encoded as **two** words.

Ok, so now you have list of googlers and projects, cruises, and sports in common. You may assume that if any two googlers have worked on the same project, gone on the same cruise, or played the same sport, they are **connected**. In other words, acquainted with each other.

Sergey Brin tells you he'll be back in a few hours to check on your work. You think to yourself *great, he's testing me* :-(. Fresh from your python class at Harvard Extension, you open the lid of your laptop, fire up a jupyter notebook, and dive in!

This is what Sergey Brin is asking you to figure out:

Part 1.

1. a) Read in the three excel files provided as panda dataframes. Hint: you might need to install the xlrd package to get excel support. List Ian Isla's projects, cruises, and sport clubs at Google. Find *one* google employee that is connected with Isla because he or she has been on the same project, cruise, or sport club. **(10 points)**
- b) Find the total number of connections at Google. We define connections as pairs of googlers that worked on a project together, had a cruise together or participated in a sports club together. For example, if Ian Isla and Kaleb Lauren went on a cruise together, then (Ian Isla, Kaleb Lauren) represents a connection. Similarly, if Kaleb Lauren and Roman Valentina played a sport together, then (Kaleb Lauren, Roman Valentina) represents a connection. Furthermore, if Ian Isla and Kaleb Lauren worked on a project together, they have another connection as well. Thus, we have another (Ian Isla, Kaleb Lauren) connection. For this problem, it might help to assign ID's to all your google employees and then loop through projects, cruises and sports to find all possible connections. **(15 points)**
- c) Create a connections dictionary where for each googler, we store a list of all connections they participate in. Using this dictionary, determine the average number of connections per googler. **(15 points)**
- d) Make use of this dictionary to find the most connected googler. What about the 10 *most* connected googlers? What about the 10 *least* connected googlers? What's the total number of connections, and the average number of connections at Google? Note that connections only counts number of projects, cruises, or sports in common. Who one is connected to makes no difference. A connection to Sergey Brin himself counts the same as a connection to a student co-op. **(10 points)**

Now, what if having important connections is of significance? In other words, being connected to Sergey Brin makes a huge difference: if a googler acquaintance you are connected to is *very* connected, that automatically makes you very **central** to google. To answer these questions, we need our old friend PageRank.

- e) One possible way to tap into PageRank is to use the now-familiar package `networkx` . We want you to create a graph from our connections dictionary. Hint: There are many ways of doing this. One possible way involves hooking into the `networkx` method `parse_adjlist` that takes in an adjacency list(connections). Look up documentation on this function to learn more. **(10 points)**
- f) Use the template from the midterm and previous assignments, or just use `networkx`'s `pagerank()` API that professor mentioned in class to calculate pageRank for google employees. Using this information, find the 10 most central googlers, and the 10 least central googlers. Are they the same as the 10 most and 10 least connected googlers? **(10 points)**

g) BONUS: Can you plot google's connectome? That connectome should easily reveal the least central googlers, since they represent isolated nodes, far away from network traffic. You have some latitude on how you choose to do this. Can you do something to upend that graph to also easily reveal the most central googlers? **(5 points)**

You meet Larry Page

While you're pondering, Larry Page comes by your desk. He says wants to see even more connected googlers, and asks you to design a googler-you-may-know recommendation engine.



You think to yourself that... just maybe googlers might be otherwise (undocumentedly) connected with connections of their connections (outside of projects, cruises, and sports clubs)! So you decide to write some python to assemble connections of connections.

Part 2.

2. a) Write a list comprehension (ideally) or some type of loop to determine connections of connections. Essentially, you should already have a list of connections, but now search the connections of a specific googler to find connections of connections. For example, if (Ian Isla, Kaleb Lauren) were connected in Part 1, and (Kalen Lauren, Roman Valentina) were also connected, then Roman Valentina would be a "connection of a connection" of Ian Isla. While you do this, perform a check to ensure that the connection of a connection should not already be a connection itself. **(20 points)**

b) What's the average number of connections-of-connections at google? You will be surprised to find out that even with only 3,000 employees, and just a few projects, cruises, and sports in common, the number of connections^{2 2} is huge. So, it would take forever to iterate through all possible employees. DON'T DO THIS. Instead, just use a random sampling technique to *sample* through a fraction of employees, throw in some debugging print statements in your loop(s) to see how fast your sampling is progressing, produce a few random samplings and average over these samplings. We leave it up to you how many samples you decide to choose, and will be lenient with our grading as long as your procedure seems reasonable. **(10 points)**

c) BONUS: What's the average distance that separates any two googlers? You may want to read [this](http://blogs.cornell.edu/info2040/2011/09/26/networks-in-hollywood-the-six-degrees-of-kevin-bacon/) (<http://blogs.cornell.edu/info2040/2011/09/26/networks-in-hollywood-the-six-degrees-of-kevin-bacon/>) article about Kevin Bacon. [This](https://en.wikipedia.org/wiki/Erd%C5%91s%E2%80%93Bacon_number) (https://en.wikipedia.org/wiki/Erd%C5%91s%E2%80%93Bacon_number) is a funny article, too. **(5 points)**

Congratulations! A successful first day at Google! You've made both Sergei and Larry happy. You're on your way to the top of the ladder! Exhausted, you escape from your desk before any other VP of something can ask anything else from you, and stop for some sushi and a little jug of warm sake. Tomorrow is new employee orientation, and you're looking forward to introduce yourself to the most connected googlers. Now *that's* taking advantage of insider information!

Solutions: Part 1

Note: This solution analyzes Question 1 and 2 by partitioning the datasets. See FULL version for solution to Question 1. Each cell here can be executed for demonstration purposes.

a) Read in the three excel files provided as panda dataframes. Hint: you might need to install the xlrd package to get excel support. List Ian Isla's projects, cruises, and sport clubs at Google. Find *one* google employee that is connected with Isla because he or she has been on the same project, cruise, or sport club. **(10 points)**

In [280]:

```
import pandas as pd

#gprojects are the projects googlers have worked on. Encoded as sentences of four words
#Notes: The first column is the name of google employees. The subsequent columns are
gprojects = pd.read_excel('gprojects.xlsx')

#gcrises are the cruises that googlers have gone on. Encoded as sentences of three words
#Notes: Each row has the name of the google employee in the first column. Then the columns are
gcrises = pd.read_excel('gcrises.xlsx')

#gsports are the sport clubs that googlers belong to. Encoded as sentences of two words
#Notes: The first column contains googler employee name, subsequent columns the sport club
gsports = pd.read_excel('gsports.xlsx')

display(gprojects.head(), gcrises.head(), gsports.head())
```

	0	1	2	3	4	5
Ian Isla	me opening of of	or emtheem intend melancholy	of the may color	may about may is	there lean Another the	NaN
Julian Isabel	brush rousing my truth:	no of my is	there music always-new to	that teacher em —after seems	dark no of seems	will talking the truth:
Roman Valentina	melancholy One teacher writing	person house cluster set	NaN	NaN	NaN	NaN

Diego Lyla	better cluster range at	and ground where The	have or the gray	NaN	NaN	NaN
Kaleb Lauren	em—after set poems warren	of urgency color only	poems intend an then	remembering One small one	brush is is of	of it it am
	0	1	2	3	4	
Ian Isla	NaN	NaN	NaN	NaN		NaN
Julian Isabel	and me door	front on the	lean words but	NaN		NaN
Roman Valentina	me Another she	NaN	NaN	NaN		NaN
Diego Lyla	gray the boulders	gray the boulders	front on the	may of and	instead idleness to	
Kaleb Lauren	your To gray	the place rooms	image her falling	then over chosen		the instead addressing

	0	1	2	3	
Ian Isla	Maybe When	of of	NaN	NaN	
Julian Isabel	NaN	NaN	NaN	NaN	
Roman Valentina	of of	NaN	NaN	NaN	
Diego Lyla	NaN	NaN	NaN	NaN	
Kaleb Lauren	am door	dawn Maybe	that to	is Another	

In [39]:

```
#Ian Isla's Google Cruises: None
gcruises.T.loc[:, ['Ian Isla']]
```

Out[39]:

	Ian Isla
0	NaN
1	NaN
2	NaN
3	NaN
4	NaN

In [38]:

```
#Ian Isla's Google Projects: Five
gprojects.T.loc[:,['Ian Isla']]
```

Out[38]:

	ian Isla
0	me opening of of
1	or emtheem intend melancholy
2	of the may color
3	may about may is
4	there lean Another the
5	NaN

In [40]:

```
#Ian Isla's Google Sport Clubs: Two
gsports.T.loc[:,['Ian Isla']]
```

Out[40]:

	ian Isla
0	Maybe When
1	of of
2	NaN
3	NaN

In [41]:

```
#Ian Isla is the first row and he shares the same first project as each of the below
gprojects.loc[gprojects[0] == 'me opening of of']

#Ian Isla is the first row and he shares the same Sport Club as many people! (e.g. 4)
gsports.loc[gsports[0] == 'Maybe When'].head(10)
```

Out[41]:

	0	1	2	3
Ian Isla	Maybe When	of of	NaN	NaN
Angel Addison	Maybe When	NaN	NaN	NaN
Jake Arianna	Maybe When	the touch	the touch	NaN
Diego Ruby	Maybe When	am door	NaN	NaN
Bradley Serenity	Maybe When	is Another	NaN	NaN
Bradley Eleanor	Maybe When	is of	dawn Maybe	over addressing
Lucas Chloe	Maybe When	towards over	NaN	NaN
Jonathan Alana	Maybe When	gray me	from the	NaN
Isaiah Eleanor	Maybe When	NaN	NaN	NaN
Cayden Adalyn	Maybe When	that to	NaN	NaN

b) Find the total number of connections at Google. We define connections as pairs of googlers that worked on a project together, had a cruise together or participated in a sports club together. For example, if Ian Isla and Kaleb Lauren went on a cruise together, then (Ian Isla, Kaleb Lauren) represents a connection. Similarly, if Kaleb Lauren and Roman Valentina played a sport together, then (Kaleb Lauren, Roman Valentina) represents a connection. Furthermore, if Ian Isla and Kaleb Lauren worked on a project together, they have another connection as well. Thus, we have another (Ian Isla, Kaleb Lauren) connection. For this problem, it might help to assign ID's to all your google employees and then loop through projects, cruises and sports to find all possible connections. **(15 points)**

In [281]:

```
#Move the current index to a new column and reindex (to create IDs), then move employee  
#Bring the Index over to create EmployeeID where every Employee has a unique value  
gsports['employee'] = gsports.index  
gsports = gsports.reset_index(drop=True)  
gsports['employeeID'] = gsports.index  
gsports = gsports[['employeeID', 'employee', 0, 1, 2, 3]]  
  
gcruises['employee'] = gcruises.index  
gcruises = gcruises.reset_index(drop=True)  
gcruises['employeeID'] = gcruises.index  
gcruises = gcruises[['employeeID', 'employee', 0, 1, 2, 3, 4]]  
  
gprojects['employee'] = gprojects.index  
gprojects = gprojects.reset_index(drop=True)  
gprojects['employeeID'] = gprojects.index  
gprojects = gprojects[['employeeID', 'employee', 0, 1, 2, 3, 4, 5]]  
  
#Merge files together on Employee  
#google = pd.merge(gprojects, pd.merge(gcruises, gsports, on = 'employee'), on = 'employeeID')  
#google
```

In [282]:

```
#Run time for this program is overwhelmingly slow... So I have truncated the data in  
gsports = gsports.truncate(before=0, after=100)  
gcruises = gcruises.truncate(before=0, after=100)  
gprojects = gprojects.truncate(before=0, after=100)
```

In [283]:

```
#Create a function to sort through and return True/False based on matches between the  
def valueMatch(val, list):  
    for col in list:  
        if col == val:  
            return True  
    return False
```

In [320]:

```
#For this problem, we will need to count the amount of matches in each column and in  
#Assumptions:  
#1) EmployeeIDs are created from the pandas index. It has been verified that the same  
#the exact same order in EACH dataset. Therefore, the same EmployeeID represents one  
#2) Matched employee pairs can count for more than one groups (so the same employee  
#can count in up to three groups - Sports, Cruises or Projects)  
#3) There is only one possible match per Employee pair per group (regardless of how  
#Cruises or Projects shared)  
  
dict = {} #Dictionary to hold name matches (Employee 1, Employee 2) and (Employee 2,  
dictID = {} #Dictionary to hold a list of matches per Employee
```



```

#Sports Club Connections
sportsCount = 0
#Iterate through all rows
for i in range(len(gsports.index)-1):
    #Compare row being iterated to the next rows
    for j in range(i+1, len(gsports.index)):
        #Transpose the first record and provide the array into the valueMatch functi
        for value in gsports.T[i].values:
            if valueMatch(value, gsports.T[j].values):
                #Add a count to the match!
                sportsCount += 1
                #We have to add the matches for both Employees, but the counter wil
                #Add in (i, j)
                try:
                    #If the key DOES exist, append the pair to the existing key val
                    dict[gsports['employee'][i]].append([gsports['employee'][i], gsports['employee'][j]])
                    dictID[gsports['employeeID'][i]].append(gsports['employeeID'][j])
                except KeyError:
                    #If a key DOES NOT exist, create it and assign the first matched
                    dict[gsports['employee'][i]] = [[gsports['employee'][i], gsports['employee'][j]]]
                    dictID[gsports['employeeID'][i]] = [gsports['employeeID'][j]]
                #Add in (j, i)
                try:
                    dict[gsports['employee'][j]].append([gsports['employee'][j], gsports['employee'][i]])
                    dictID[gsports['employeeID'][j]].append(gsports['employeeID'][i])
                except KeyError:
                    dict[gsports['employee'][j]] = [[gsports['employee'][j], gsports['employee'][i]]]
                    dictID[gsports['employeeID'][j]] = [gsports['employeeID'][i]]
                #If there is any match at all, these employees are a matched connect
                break

```

```

#Cruise Connections (Follows same logic as above)
cruisesCount = 0
for i in range(len(gcruises.index)-1):
    for j in range(i+1, len(gcruises.index)):
        for value in gcruises.T[i].values:
            if valueMatch(value, gcruises.T[j].values):
                cruisesCount += 1
                try:
                    dict[gcruises['employee'][i]].append([gcruises['employee'][i], gcruises['employee'][j]])
                    dictID[gcruises['employeeID'][i]].append(gcruises['employeeID'][j])
                except KeyError:
                    dict[gcruises['employee'][i]] = [[gcruises['employee'][i], gcruises['employee'][j]]]
                    dictID[gcruises['employeeID'][i]] = [gcruises['employeeID'][j]]
                try:
                    dict[gcruises['employee'][j]].append([gcruises['employee'][j], gcruises['employee'][i]])
                    dictID[gcruises['employeeID'][j]].append(gcruises['employeeID'][i])
                except KeyError:
                    dict[gcruises['employee'][j]] = [[gcruises['employee'][j], gcruises['employee'][i]]]
                    dictID[gcruises['employeeID'][j]] = [gcruises['employeeID'][i]]
                break

```

```

#Project Connections (Follows same logic as above)
projectsCount = 0

```

```

for i in range(len(gprojects.index)-1):
    for j in range(i+1, len(gprojects.index)):
        for value in gprojects.T[i].values:
            if valueMatch(value, gprojects.T[j].values):
                projectsCount += 1
                try:
                    dict[gprojects['employee'][i]].append([gprojects['employee'][i], gprojects['employeeID'][i]])
                    dictID[gprojects['employeeID'][i]].append(gprojects['employeeID'][j])
                except KeyError:
                    dict[gprojects['employee'][i]] = [[gprojects['employee'][i], gprojects['employeeID'][i]]]
                    dictID[gprojects['employeeID'][i]] = [gprojects['employeeID'][j]]
                try:
                    dict[gprojects['employee'][j]].append([gprojects['employee'][j], gprojects['employeeID'][j]])
                    dictID[gprojects['employeeID'][j]].append(gprojects['employeeID'][i])
                except KeyError:
                    dict[gprojects['employee'][j]] = [[gprojects['employee'][j], gprojects['employeeID'][j]]]
                    dictID[gprojects['employeeID'][j]] = [gprojects['employeeID'][i]]
            break

print("Count of Connections from Sports Clubs: "+str(sportsCount))
print("Count of Connections from Cruises: "+str(cruisesCount))
print("Count of Connections from Projects: "+str(projectsCount))

```

```

Count of Connections from Sports Clubs: 622
Count of Connections from Cruises: 230
Count of Connections from Projects: 90

```

c) Create a connections dictionary where for each googler, we store a list of all connections they participate in. Using this dictionary, determine the average number of connections per googler. **(15 points)**

In [108]:

```
#Dictionary created as a part of Part B
dict
```

Out[108]:

```
{'Ian Isla': [['Ian Isla', 'Roman Valentina'],
              ['Ian Isla', 'Christian Alaina'],
              ['Ian Isla', 'Emmett Aurora'],
              ['Ian Isla', 'Angel Addison'],
              ['Ian Isla', 'Jake Arianna'],
              ['Ian Isla', 'George Bella'],
              ['Ian Isla', 'Hudson Melanie'],
              ['Ian Isla', 'Emmett Riley'],
              ['Ian Isla', 'Wesley Athena'],
              ['Ian Isla', 'Ashton Quinn'],
              ['Ian Isla', 'Luca Adriana'],
              ['Ian Isla', 'Aiden Sarah']],
 'Roman Valentina': [['Roman Valentina', 'Emmett Aurora'],
                     ['Roman Valentina', 'Hudson Melanie'],
                     ['Roman Valentina', 'Emmett Riley'],
                     ['Roman Valentina', 'Wesley Athena'],
                     ['Roman Valentina', 'Ashton Quinn'],
                     ['Roman Valentina', 'Timothy McKenzie']]}
```

In [326]:

```
total = sportsCount + cruisesCount + projectsCount
count = len(dict.keys())
average = total/count
print("Average Connections: "+str(average))
```

Average Connections: 9.42

d) Make use of this dictionary to find the most connected googler. What about the 10 *most* connected googlers? What about the 10 *least* connected googlers? What's the total number of connections, and the average number of connections at Google? Note that connections only counts number of projects, cruises, or sports in common. Who one is connected to makes no difference. A connection to Sergey Brin himself counts the same as a connection to a student co-op. **(10 points)**

In [327]:

```
print("Total Connections: "+str(total))
print("Average Connections: "+str(average))
```

Total Connections: 942

Average Connections: 9.42

In [297]:

```
#Use bubble sort from our lectures, but modify it to work with a list of tuples so  
#we can order key-value pairs from greatest to least  
def bubble_sort(alist):  
    for passnum in range(len(alist)-1, -1, -1):  
        swapped = False  
        for i in range(passnum):  
            if alist[i][1] < alist[i+1][1]:  
                alist[i], alist[i+1] = alist[i+1], alist[i]  
                swapped = True  
        if not swapped:  
            break
```

In [298]:

```
#Pass the key AND # of contents into a list of tuples  
sorted_list = []  
for key, value in sorted(dict.items()):  
    sorted_list.append((key, len([item for item in value if item])))  
bubble_sort(sorted_list)
```

In [328]:

```
#Top Ten Connected Googlers!  
sorted_list[:10]
```

Out[328]:

```
[('Joel Ava', 52),  
 ('Owen Adriana', 46),  
 ('Jameson Trinity', 45),  
 ('Colton Ava', 43),  
 ('Logan Katelyn', 43),  
 ('John Eliza', 42),  
 ('Isaac Makenzie', 41),  
 ('Timothy Jayla', 39),  
 ('Preston Adalynn', 38),  
 ('Alan Ella', 36)]
```

In [300]:

```
#Bottom Ten Connected Googlers!  
sorted_list[-10:]
```

Out[300]:

```
[('Nathan Savannah', 6),  
 ('Colton Valentina', 5),  
 ('Weston Ariel', 5),  
 ('Cayden Quinn', 4),  
 ('Jayce Bella', 3),  
 ('Micah Morgan', 2),  
 ('Thomas Peyton', 2),  
 ('Damian Peyton', 1),  
 ('George Ruby', 1),  
 ('Leonardo Valeria', 1)]
```

Now, what if having important connections is of significance? In other words, being connected to Sergey Brin makes a huge difference: if a googler acquaintance you are connected to is *very* connected, that automatically makes you very ***central*** to google. To answer these questions, we need our old friend PageRank.

e) One possible way to tap into PageRank is to use the now-familiar package `networkx`. We want you to create a graph from our connections dictionary. Hint: There are many ways of doing this. One possible way involves hooking into the `networkx` method `parse_adjlist` that takes in an adjacency list(connections). Look up documentation on this function to learn more. **(10 points)**

In [336]:

```
#Dictionary of EmployeeIDs created as a part of Part B  
#Note: This is due to parse_adjust needing numbers to split nodes/edges  
#Names ('First Last') would be two entities 'First' as a node and 'Last' as an edge  
print(dictID)
```

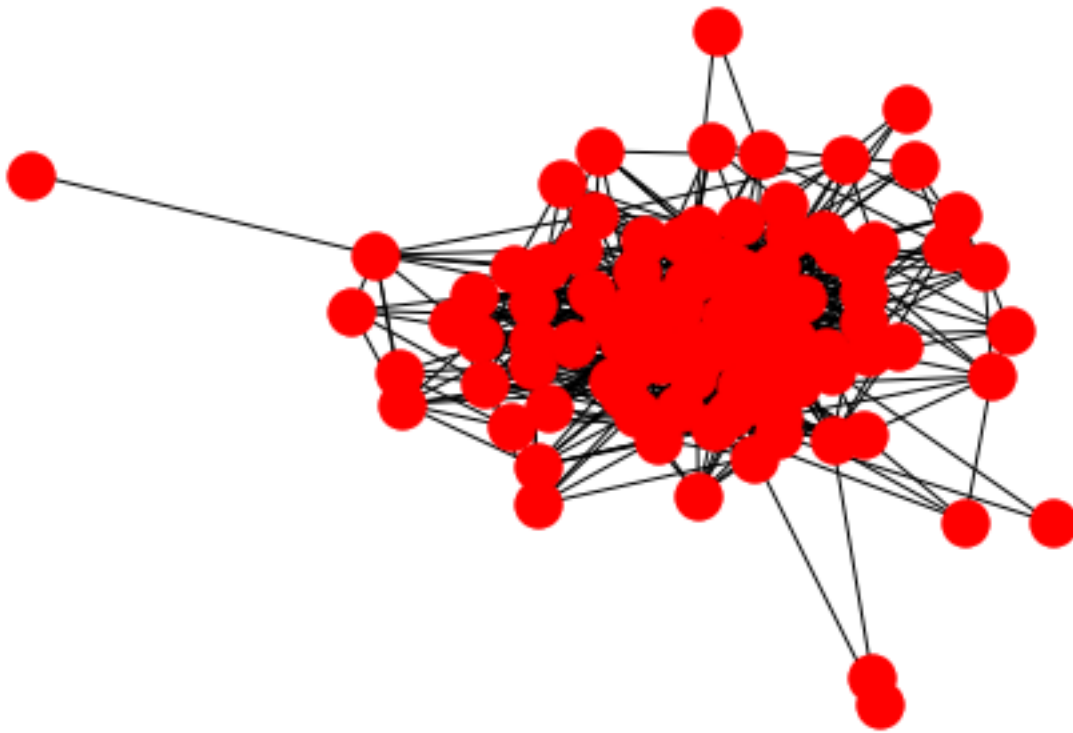
```
{0: [2, 6, 9, 26, 27, 58, 65, 81, 83, 93, 62, 70], 2: [0, 9, 65, 81, 83, 93, 21], 6: [0, 16, 21, 23, 24, 25, 26, 27, 29, 39, 42, 48, 50, 51, 58, 61, 62, 67, 68, 81, 18, 37, 73], 9: [0, 2, 8, 14, 18, 25, 65, 67, 69, 78, 81, 83, 88, 92, 93, 21, 27, 28, 50, 77, 11, 63, 77, 85], 26: [0, 6, 27, 58, 35, 76, 32, 37, 50], 27: [0, 6, 11, 26, 30, 36, 58, 66, 79, 83, 87, 92, 93, 4, 9, 35, 50, 54, 94, 96, 92], 58: [0, 6, 10, 11, 17, 21, 23, 26, 27, 30, 36, 46, 48, 50, 62, 66, 70, 79, 83, 87, 92, 93, 29, 65, 99, 66], 65: [0, 2, 4, 5, 9, 10, 16, 19, 25, 38, 39, 52, 60, 71, 81, 83, 86, 87, 90, 93, 95, 100, 29, 40, 46, 58, 77, 24, 46], 81: [0, 2, 6, 9, 21, 23, 24, 25, 39, 42, 48, 50, 51, 61, 62, 65, 67, 83, 93, 40], 83: [0, 2, 9, 11, 27, 30, 36, 58, 65, 66, 79, 81, 87, 92, 93, 1, 35, 35], 93: [0, 2, 9, 11, 27, 30, 36, 58, 60, 61, 65, 66, 75, 79, 81, 83, 87, 92, 12], 4: [5, 10, 11, 12, 16, 19, 25, 30, 38, 39, 52, 54, 59, 60, 65, 70, 71, 85, 86, 87, 90, 92, 95, 97, 100, 14, 20, 27, 35, 38, 51, 55, 68, 71], 5: [4, 10, 11, 12, 16, 19, 25, 30, 38, 39, 40, 42, 50, 52, 56, 60, 65, 70, 71, 84, 86, 87, 90, 92, 95, 100, 39, 46, 51, 94], 10: [4, 5, 14, 16, 17, 19, 21, 23, 25, 30, 38, 39, 46, 48, 50, 52, 58, 60, 62, 65, 70, 71, 86, 87, 90, 92, 95, 100], 11: [4, 5, 12, 27, 30, 36, 39, 58, 66, 70, 79, 83, 87, 92, 93, 3, 31, 33, 54, 74, 78, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100]}
```

In [302]:

```
import networkx as nx  
  
#Append the employeeIDs to create a list of strings (e.g. ['1 2 3', '2 2']) with key  
test = []  
for key, value in sorted(dictID.items()):  
    test.append(str(key)+" "+" ".join([str(item) for item in value if item]))  
  
#https://networkx.github.io/documentation/networkx-1.10/reference/generated/networkx  
ggraph = nx.parse_adjlist(test, nodetype=str)
```

In [303]:

```
#Graph of Connections Drawn Out!  
nx.draw(ggraph)
```



f) Use the template from the midterm and previous assignments, or just use networkx's `pagerank()` API that professor mentioned in class to calculate pageRank for google employees. Using this information, find the 10 most central googlers, and the 10 least central googlers. Are they the same as the 10 most and 10 least connected googlers? **(10 points)**

In [304]:

```
import numpy as np  
  
#Apply PageRank to the graph of Googlers  
ggraph_pr = nx.pagerank(ggraph) #Note: Default damping is 0.85!  
  
#Convert the dictionary to a tuple and sort  
ggraph_ranked = [(k, v) for k, v in ggraph_pr.items()]  
bubble_sort(ggraph_ranked)
```

In [305]:

```
#Redo Part D to sort with EmployeeID rather than Employee names  
empID_connections = []  
for key, value in sorted(dictID.items()):  
    empID_connections.append((key, len([item for item in value if item])))  
bubble_sort(empID_connections)
```


In [306]:

```
#Top Ten Googlers by PageRank!  
#Reference Pandas Table for Full Names associated with EmployeeIDs  
ggraph_ranked[:10]
```

Out[306]:

```
[('50', 0.022820129402410672),  
 ('60', 0.02194001836738658),  
 ('52', 0.021460032170229196),  
 ('67', 0.020625117069219122),  
 ('25', 0.02033308102254799),  
 ('16', 0.020250826019185545),  
 ('39', 0.018278153501639652),  
 ('30', 0.018227946149985987),  
 ('71', 0.016717148421104666),  
 ('4', 0.01644505221227893)]
```

In [307]:

```
#Top Ten Connected Googlers (but by EmployeeID, # of Connections)  
#Reference Pandas Table for Full Names associated with EmployeeIDs  
empID_connections[:10]
```

Out[307]:

```
[(50, 52),  
 (16, 46),  
 (52, 45),  
 (25, 43),  
 (39, 43),  
 (60, 42),  
 (67, 41),  
 (30, 39),  
 (71, 38),  
 (48, 36)]
```

In [308]:

```
#Bottom Ten Googlers by PageRank!  
#Reference Pandas Table for Full Names associated with EmployeeIDs  
ggraph_ranked[-10:]
```

Out[308]:

```
[ ('99', 0.004468905019286139),  
  ('91', 0.004172302704734399),  
  ('94', 0.003822089001090339),  
  ('64', 0.0035784682084692614),  
  ('57', 0.002957342386297933),  
  ('80', 0.002620839816228804),  
  ('15', 0.002409624501967589),  
  ('53', 0.002226273708244273),  
  ('44', 0.002080589913150963),  
  ('72', 0.001996646491383476)]
```

In [309]:

```
#Bottom Ten Connected Googlers (but by EmployeeID, # of Connections)  
#Reference Pandas Table for Full Names associated with EmployeeIDs  
empID_connections[-10:]
```

Out[309]:

```
[ (97, 6),  
  (91, 5),  
  (94, 5),  
  (64, 4),  
  (57, 3),  
  (15, 2),  
  (80, 2),  
  (44, 1),  
  (53, 1),  
  (72, 1)]
```

Note: We can see that the top and bottom pagerank and connected employee list are different (as we have noticed in class before) in that the top members of pagerank are not the same as the top members with the most connections. Those with connections indirectly others with MORE connections often increases the pagerank more - although there is much overlap between the top ten groups. Given the analogy of Sergey Brin previously - a connection to Sergey Brin at this point raises a person's pagerank more than a connection to a co-op student (assuming the co-op student has very few connections in turn) so it is more of a measure of importance!

g) BONUS: Can you plot google's connectome? That connectome should easily reveal the least central googlers, since they represent isolated nodes, far away from network traffic. You have some latitude on how you choose to do this. Can you do something to upend that graph to also easily reveal the most central googlers? **(5 points)**

```
In [ ]:
```

Solutions: Part 2

a) Write a list comprehension (ideally) or some type of loop to determine connections of connections. Essentially, you should already have a list of connections, but now search the connections of a specific googler to find connections of connections. For example, if (Ian Isla, Kaleb Lauren) were connected in Part 1, and (Kalen Lauren, Roman Valentina) were also connected, then Roman Valentina would be a "connection of a connection" of Ian Isla. While you do this, perform a check to ensure that the connection of a connection should not already be a connection itself. **(20 points)**

```
In [310]:

seconds = pd.DataFrame.from_dict(dictID, orient='index')
seconds.head()
```

Out[310]:

	0	1	2	3	4	5	6	7	8	9	...	42	43	44	45	46	47
0	2	6.0	9.0	26.0	27.0	58.0	65.0	81.0	83.0	93.0	...	NaN	NaN	NaN	NaN	NaN	NaN
2	0	9.0	65.0	81.0	83.0	93.0	21.0	NaN	NaN	NaN	...	NaN	NaN	NaN	NaN	NaN	NaN
6	0	16.0	21.0	23.0	24.0	25.0	26.0	27.0	29.0	39.0	...	NaN	NaN	NaN	NaN	NaN	NaN
9	0	2.0	8.0	14.0	18.0	25.0	65.0	67.0	69.0	78.0	...	NaN	NaN	NaN	NaN	NaN	NaN
26	0	6.0	27.0	58.0	35.0	76.0	32.0	37.0	50.0	NaN	...	NaN	NaN	NaN	NaN	NaN	NaN

5 rows × 52 columns

```
In [311]:

seconds['employeeID'] = seconds.index
seconds = seconds.reset_index(drop=True)
```

In [312]:

```
seconds.head()
```

Out[312]:

	0	1	2	3	4	5	6	7	8	9	...	43	44	45	46	47	48	
0	2	6.0	9.0	26.0	27.0	58.0	65.0	81.0	83.0	93.0	...	NaN	NaN	NaN	NaN	NaN	NaN	
1	0	9.0	65.0	81.0	83.0	93.0	21.0	NaN	NaN	NaN	...	NaN	NaN	NaN	NaN	NaN	NaN	
2	0	16.0	21.0	23.0	24.0	25.0	26.0	27.0	29.0	39.0	...	NaN	NaN	NaN	NaN	NaN	NaN	
3	0	2.0	8.0	14.0	18.0	25.0	65.0	67.0	69.0	78.0	...	NaN	NaN	NaN	NaN	NaN	NaN	
4	0	6.0	27.0	58.0	35.0	76.0	32.0	37.0	50.0	NaN	...	NaN	NaN	NaN	NaN	NaN	NaN	

5 rows × 53 columns

In [339]:

```
#Create dictionaries to hold the values
secondConn = {}

#Secondary Connections - Uses same logic as before (except for two commented lines)
secondCount = 0
for i in range(len(seconds.index)-1):
    for j in range(i+1, len(seconds.index)):
        for value in seconds.T[i].values:
            #Check that this employee is not already a Primary Connection
            if valueMatch(seconds['employeeID'][i], seconds.T[j].values) is False:
                #Check to see if there are any matching secondary connections
                if valueMatch(value, seconds.T[j].values):
                    secondCount += 1
                    try:
                        secondConn[seconds['employeeID'][i]].append(seconds['employeeID'][j])
                    except KeyError:
                        secondConn[seconds['employeeID'][i]] = [seconds['employeeID'][j]]
                    try:
                        secondConn[seconds['employeeID'][j]].append(seconds['employeeID'][i])
                    except KeyError:
                        secondConn[seconds['employeeID'][j]] = [seconds['employeeID'][i]]
                    break

print("Count of Secondary Connections: "+str(secondCount))
```

Count of Secondary Connections: 3214

In [334]:

```
print(secondConn)
```

```
{0: [4, 5, 10, 11, 12, 16, 19, 25, 30, 38, 39, 52, 54, 60, 71, 85, 86, 87, 90, 92, 95, 100, 40, 42, 50, 84, 21, 23, 24, 29, 48, 51, 61, 67, 68, 7, 8, 14, 28, 69, 18, 78, 88, 17, 46, 36, 66, 79, 75, 1, 35, 37, 73, 77, 96, 22, 64, 99, 32, 63, 76, 94], 2: [6, 26, 27, 58, 4, 5, 10, 11, 12, 16, 19, 25, 30, 38, 39, 52, 60, 70, 71, 85, 86, 87, 90, 92, 95, 97, 100, 40, 42, 50, 23, 24, 29, 48, 51, 61, 62, 67, 8, 14, 28, 69, 18, 78, 88, 17, 46, 36, 66, 79, 75, 1, 35, 77, 72, 63], 6: [9, 65, 83, 93, 3, 4, 5, 10, 11, 12, 19, 30, 38, 52, 54, 60, 70, 71, 85, 86, 87, 90, 92, 95, 97, 100, 40, 56, 84, 7, 8, 14, 20, 28, 31, 55, 69, 78, 88, 17, 46, 36, 66, 79, 13, 41, 34, 75, 1, 3, 74, 35, 77, 96, 72, 22, 43, 45, 64, 99, 32, 63, 76, 94, 89, 47, 49, 98, 57], 9: [26, 58, 4, 5, 10, 12, 16, 19, 30, 38, 39, 52, 54, 59, 60, 70, 71, 86, 87, 90, 95, 97, 100, 40, 42, 56, 84, 23, 24, 29, 48, 51, 61, 62, 68, 7, 20, 31, 55, 17, 46, 36, 66, 79, 13, 34, 75, 1, 3, 33, 74, 35, 37, 73, 96, 72, 22, 43, 45, 64, 99, 76, 94, 89, 47, 98, 44], 26: [65, 81, 83, 93, 4, 5, 10, 11, 12, 16, 25, 30, 39, 52, 54, 60, 70, 71, 85, 87, 92, 40, 42, 56, 84, 21, 23, 24, 29, 48, 51, 61, 62, 67, 68, 7, 8, 14, 20, 31, 55, 69, 18, 78, 17, 46, 36, 66, 79, 75, 1, 74, 73, 96, 99, 63, 94, 89, 98, 53, 91], 27: [65, 81, 5, 10, 12, 16, 19, 25, 38, 39, 52, 59, 60, 70, 71, 85, 86, 87, 90, 92, 95, 97, 100, 40, 42, 56, 84, 21, 23, 24, 29, 48, 51, 61, 62, 67, 68, 7, 8, 14, 20, 31, 55, 69, 18, 78, 17, 46, 36, 66, 79, 75, 1, 74, 73, 96, 99, 63, 94, 89, 98, 53, 91]}
```

b) What's the average number of connections-of-connections at google? You will be surprised to find out that even with only 3,000 employees, and just a few projects, cruises, and sports in common, the number of connections²² is huge. So, it would take forever to iterate through all possible employees. DON'T DO THIS. Instead, just use a random sampling technique to *sample* through a fraction of employees, throw in some debugging print statements in your loop(s) to see how fast your sampling is progressing, produce a few random samplings and average over these samplings. We leave it up to you how many samples you decide to choose, and will be lenient with our grading as long as your procedure seems reasonable. **(10 points)**

In [340]:

```
#Since I am only beginning with a small subset of the original dataset, it is much easier to sample
total = secondCount

#Not everyone would presumably have secondary connections, so we use the original array to count all
countAll = len(dict.items())
averageAll = total/countAll

#Only out of the pool of people with secondary connections
countSeconds = len(secondConn.items())
averageSeconds = total/countSeconds

print("Average Secondary Connections (Amongst Everyone): "+str(averageAll))
print("Average Secondary Connections (Only Amongst Googlers With Secondary Connections): "+str(averageSeconds))
```

```
Average Secondary Connections (Amongst Everyone): 32.14
Average Secondary Connections (Only Amongst Googlers With Secondary Connections): 32.14
```

Note: in this case everyone has a secondary connection - but that is not necessarily true!

c) BONUS: What's the average distance that separates any two googlers? You may want to read [this](http://blogs.cornell.edu/info2040/2011/09/26/networks-in-hollywood-the-six-degrees-of-kevin-bacon/) (<http://blogs.cornell.edu/info2040/2011/09/26/networks-in-hollywood-the-six-degrees-of-kevin-bacon/>) article about Kevin Bacon. [This](https://en.wikipedia.org/wiki/Erd%C5%91s%E2%80%93Bacon_number) (https://en.wikipedia.org/wiki/Erd%C5%91s%E2%80%93Bacon_number) is a funny article, too. **(5 points)**

In []:

CSCI E-7 Introduction to Python Programming for Life Sciences
Dino Konstantopoulos, **Final Exam, Take Home Portion**

Working at Google - good idea?

You meet Sergey Brin

You've been hired at Google to run their **Life Sciences division**, congratulations! On your first day after your move to California's Silicon Valley, **Sergey Brin** comes in to meet the new wizard pythonista. He gets talking and admits to you there's so many googlers that he just cannot keep track of all of them. But he would like to build a rolodex of **central** googlers. Not just central in terms of titles or degrees or awards, but central in terms of social connections. Sergey wants you to find out who the most *social* googlers are.



Sergey has an idea: He gives you a list of projects that googlers have worked on, a list of Cruises googlers have gone on together, and a list of sport clubs users belong to.

You get 3 excel files:

- 1) `gprojects` are the projects googlers have worked on. The first column is the name of google employees, the subsequent columns are the projects that they have worked on. Googlers love codes, so the projects are encoded as sentences of **four** words each. The sentences don't need to make sense.
- 2) `gcruises` are the cruises that googlers have gone on. each row has the name of the google employee in the first column, and then the cruises she has been on. Cruises are encoded as sentences of **three** words each.
- 3) `gsports` are the sport clubs that googlers belong to. The first column contains googler employee name, subsequent columns the sport clubs they belong to. Sport clubs are encoded as **two** words.

Ok, so now you have list of googlers and projects, cruises, and sports in common. You may assume that if any two googlers have worked on the same project, gone on the same cruise, or played the same sport, they are **connected**. In other words, acquainted with each other.

Sergey Brin tells you he'll be back in a few hours to check on your work. You think to yourself *great, he's testing me* :-(. Fresh from your python class at Harvard Extension, you open the lid of your laptop, fire up a jupyter notebook, and dive in!

This is what Sergey Brin is asking you to figure out:

Part 1.

1. a) Read in the three excel files provided as panda dataframes. Hint: you might need to install the `xlrd` package to get excel support. List Ian Isla's projects, cruises, and sport clubs at Google. Find *one* google employee that is connected with Isla because he or she has been on the same project, cruise, or sport club. **(10 points)**
- b) Find the total number of connections at Google. We define connections as pairs of googlers that worked on a project together, had a cruise together or participated in a sports club together. For example, if Ian Isla and Kaleb Lauren went on a cruise together, then (Ian Isla, Kaleb Lauren) represents a connection. Similarly, if Kaleb Lauren and Roman Valentina played a sport together, then (Kaleb Lauren, Roman Valentina) represents a connection. Furthermore, if Ian Isla and Kaleb Lauren worked on a project together, they have another connection as well. Thus, we have another (Ian Isla, Kaleb Lauren) connection. For this problem, it might help to assign ID's to all your google employees and then loop through projects, cruises and sports to find all possible connections. **(15 points)**
- c) Create a connections dictionary where for each googler, we store a list of all connections they participate in. Using this dictionary, determine the average number of connections per googler. **(15 points)**
- d) Make use of this dictionary to find the most connected googler. What about the 10 *most* connected googlers? What about the 10 *least* connected googlers? What's the total number of connections, and the average number of connections at Google? Note that connections only counts number of projects, cruises, or sports in common. Who one is connected to makes no difference. A connection to Sergey Brin himself counts the same as a connection to a student co-op. **(10 points)**

Now, what if having important connections is of significance? In other words, being connected to Sergey Brin makes a huge difference: if a googler acquaintance you are connected to is *very* connected, that automatically makes you very **central** to google. To answer these questions, we need our old friend PageRank.

e) One possible way to tap into PageRank is to use the now-familiar package `networkx`. We want you to create a graph from our connections dictionary. Hint: There are many ways of doing this. One possible way involves hooking into the `networkx` method `parse_adjlist` that takes in an adjacency list(connections). Look up documentation on this function to learn more. **(10 points)**

f) Use the template from the midterm and previous assignments, or just use `networkx`'s `pagerank()` API that professor mentioned in class to calculate pageRank for google employees. Using this information, find the 10 most central googlers, and the 10 least central googlers. Are they the same as the 10 most and 10 least connected googlers? **(10 points)**

g) BONUS: Can you plot google's connectome? That connectome should easily reveal the least central googlers, since they represent isolated nodes, far away from network traffic. You have some latitude on how you choose to do this. Can you do something to upend that graph to also easily reveal the most central googlers? **(5 points)**

You meet Larry Page

While you're pondering, Larry Page comes by your desk. He says wants to see even more connected googlers, and asks you to design a googler-you-may-know recommendation engine.



You think to yourself that... just maybe googlers might be otherwise (undocumentedly) connected with connections of their connections (outside of projects, cruises, and sports clubs)! So you decide to write some python to assemble connections of connections.

Part 2.

2. a) Write a list comprehension (ideally) or some type of loop to determine connections of connections. Essentially, you should already have a list of connections, but now search the connections of a specific googler to find connections of connections. For example, if (Ian Isla, Kaleb Lauren) were connected in Part 1, and (Kalen Lauren, Roman Valentina) were also connected, then Roman Valentina would be a "connection of a connection" of Ian Isla. While you do this, perform a check to ensure that the connection of a connection should not already be a connection itself. **(20 points)**
- b) What's the average number of connections-of-connections at google? You will be surprised to find out that even with only 3,000 employees, and just a few projects, cruises, and sports in common, the number of connections² is huge. So, it would take forever to iterate through all possible employees. DON'T DO THIS. Instead, just use a random sampling technique to *sample* through a fraction of employees, throw in some debugging print statements in your loop(s) to see how fast your sampling is progressing, produce a few random samplings and average over these samplings. We leave it up to you how many samples you decide to choose, and will be lenient with our grading as long as your procedure seems reasonable. **(10 points)**
- c) BONUS: What's the average distance that separates any two googlers? You may want to read [this](http://blogs.cornell.edu/info2040/2011/09/26/networks-in-hollywood-the-six-degrees-of-kevin-bacon/) (<http://blogs.cornell.edu/info2040/2011/09/26/networks-in-hollywood-the-six-degrees-of-kevin-bacon/>) article about Kevin Bacon. [This](https://en.wikipedia.org/wiki/Erd%C5%91s%E2%80%93Bacon_number) (https://en.wikipedia.org/wiki/Erd%C5%91s%E2%80%93Bacon_number) is a funny article, too. **(5 points)**

Congratulations! A successful first day at Google! You've made both Sergei and Larry happy. You're on your way to the top of the ladder! Exhausted, you escape from your desk before any other VP of something can ask anything else from you, and stop for some sushi and a little jug of warm sake. Tomorrow is new employee orientation, and you're looking forward to introduce yourself to the most connected googlers. Now *that's* taking advantage of insider information!

Solutions: Part 1

Note: This solution analyzes Question 1 and 2 by partitioning the datasets. See FULL version for solution to Question 1. Each cell here can be executed for demonstration purposes.

- a) Read in the three excel files provided as panda dataframes. Hint: you might need to install the xlrd package to get excel support. List Ian Isla's projects, cruises, and sport clubs at Google. Find *one* google employee that is connected with Isla because he or she has been on the same project, cruise, or sport club. **(10 points)**

	0	1	2	3	4	5
Ian Isla	me opening of of	or emtheem intend melancholy	of the may color	may about may is	there lean Another the	NaN
Julian Isabel	brush rousing my truth:	no of my is	there music always-new to	that teacher em —after seems	dark no of seems	will talking the truth:
Roman Valentina	melancholy One teacher writing	person house cluster set	NaN	NaN	NaN	NaN
Diego Lyla	better cluster range at	and ground where The	have or the gray	NaN	NaN	NaN
Kaleb Lauren	em—after set poems warren	of urgency color only	poems intend an then	remembering One small one	brush is is of	of it it am

	0	1	2	3	4
Ian Isla	NaN	NaN	NaN	NaN	NaN
Julian Isabel	and me door	front on the	lean words but	NaN	NaN
Roman Valentina	me Another she	NaN	NaN	NaN	NaN
Diego Lyla	gray the boulders	gray the boulders	front on the	may of and	instead idleness to
Kaleb Lauren	your To gray	the place rooms	image her falling	then over chosen	the instead addressing

	0	1	2	3
Ian Isla	Maybe When	of of	NaN	NaN
Julian Isabel	NaN	NaN	NaN	NaN
Roman Valentina	of of	NaN	NaN	NaN
Diego Lyla	NaN	NaN	NaN	NaN
Kaleb Lauren	am door	dawn Maybe	that to	is Another

In [39]:

Out[39]:

lan Isla	
0	NaN
1	NaN
2	NaN
3	NaN
4	NaN

In [38]:

Out[38]:

lan Isla	
0	me opening of of
1	or emtheem intend melancholy
2	of the may color
3	may about may is
4	there lean Another the
5	NaN

In [40]:

Out[40]:

lan Isla	
0	Maybe When
1	of of
2	NaN
3	NaN

In [41]:

Out[41]:

	0	1	2	3
Ian Isla	Maybe When	of of	NaN	NaN
Angel Addison	Maybe When	NaN	NaN	NaN
Jake Arianna	Maybe When	the touch	the touch	NaN
Diego Ruby	Maybe When	am door	NaN	NaN
Bradley Serenity	Maybe When	is Another	NaN	NaN
Bradley Eleanor	Maybe When	is of	dawn Maybe	over addressing
Lucas Chloe	Maybe When	towards over	NaN	NaN
Jonathan Alana	Maybe When	gray me	from the	NaN
Isaiah Eleanor	Maybe When	NaN	NaN	NaN
Cayden Adalyn	Maybe When	that to	NaN	NaN

b) Find the total number of connections at Google. We define connections as pairs of googlers that worked on a project together, had a cruise together or participated in a sports club together. For example, if Ian Isla and Kaleb Lauren went on a cruise together, then (Ian Isla, Kaleb Lauren) represents a connection. Similarly, if Kaleb Lauren and Roman Valentina played a sport together, then (Kaleb Lauren, Roman Valentina) represents a connection. Furthermore, if Ian Isla and Kaleb Lauren worked on a project together, they have another connection as well. Thus, we have another (Ian Isla, Kaleb Lauren) connection. For this problem, it might help to assign ID's to all your google employees and then loop through projects, cruises and sports to find all possible connections. **(15 points)**

In [281]:

In [282]:

In [283]:

In [320]:

```
Count of Connections from Sports Clubs: 622
Count of Connections from Cruises: 230
Count of Connections from Projects: 90
```

c) Create a connections dictionary where for each googler, we store a list of all connections they participate in. Using this dictionary, determine the average number of connections per googler. **(15 points)**

In [108]:

Out[108]:

```
{'Ian Isla': [['Ian Isla', 'Roman Valentina'],
              ['Ian Isla', 'Christian Alaina'],
              ['Ian Isla', 'Emmett Aurora'],
              ['Ian Isla', 'Angel Addison'],
              ['Ian Isla', 'Jake Arianna'],
              ['Ian Isla', 'George Bella'],
              ['Ian Isla', 'Hudson Melanie'],
              ['Ian Isla', 'Emmett Riley'],
              ['Ian Isla', 'Wesley Athena'],
              ['Ian Isla', 'Ashton Quinn'],
              ['Ian Isla', 'Luca Adriana'],
              ['Ian Isla', 'Aiden Sarah']],
 'Roman Valentina': [['Roman Valentina', 'Emmett Aurora'],
                     ['Roman Valentina', 'Hudson Melanie'],
                     ['Roman Valentina', 'Emmett Riley'],
                     ['Roman Valentina', 'Wesley Athena'],
                     ['Roman Valentina', 'Ashton Quinn'],
                     ['Roman Valentina', 'Timothy McKenzie']]}
```

In [326]:

Average Connections: 9.42

d) Make use of this dictionary to find the most connected googler. What about the 10 *most* connected googlers? What about the 10 *least* connected googlers? What's the total number of connections, and the average number of connections at Google? Note that connections only counts number of projects, cruises, or sports in common. Who one is connected to makes no difference. A connection to Sergey Brin himself counts the same as a connection to a student co-op. **(10 points)**

In [327]:

Total Connections: 942
Average Connections: 9.42

In [297]:

In [298]:

In [328]:

Out[328]:

```
[('Joel Ava', 52),  
 ('Owen Adriana', 46),  
 ('Jameson Trinity', 45),  
 ('Colton Ava', 43),  
 ('Logan Katelyn', 43),  
 ('John Eliza', 42),  
 ('Isaac Makenzie', 41),  
 ('Timothy Jayla', 39),  
 ('Preston Adalynn', 38),  
 ('Alan Ella', 36)]
```


In [300]:

Out[300]:

```
[('Nathan Savannah', 6),  
 ('Colton Valentina', 5),  
 ('Weston Ariel', 5),  
 ('Cayden Quinn', 4),  
 ('Jayce Bella', 3),  
 ('Micah Morgan', 2),  
 ('Thomas Peyton', 2),  
 ('Damian Peyton', 1),  
 ('George Ruby', 1),  
 ('Leonardo Valeria', 1)]
```

Now, what if having important connections is of significance? In other words, being connected to Sergey Brin makes a huge difference: if a googler acquaintance you are connected to is *very* connected, that automatically makes you very ***central*** to google. To answer these questions, we need our old friend PageRank.

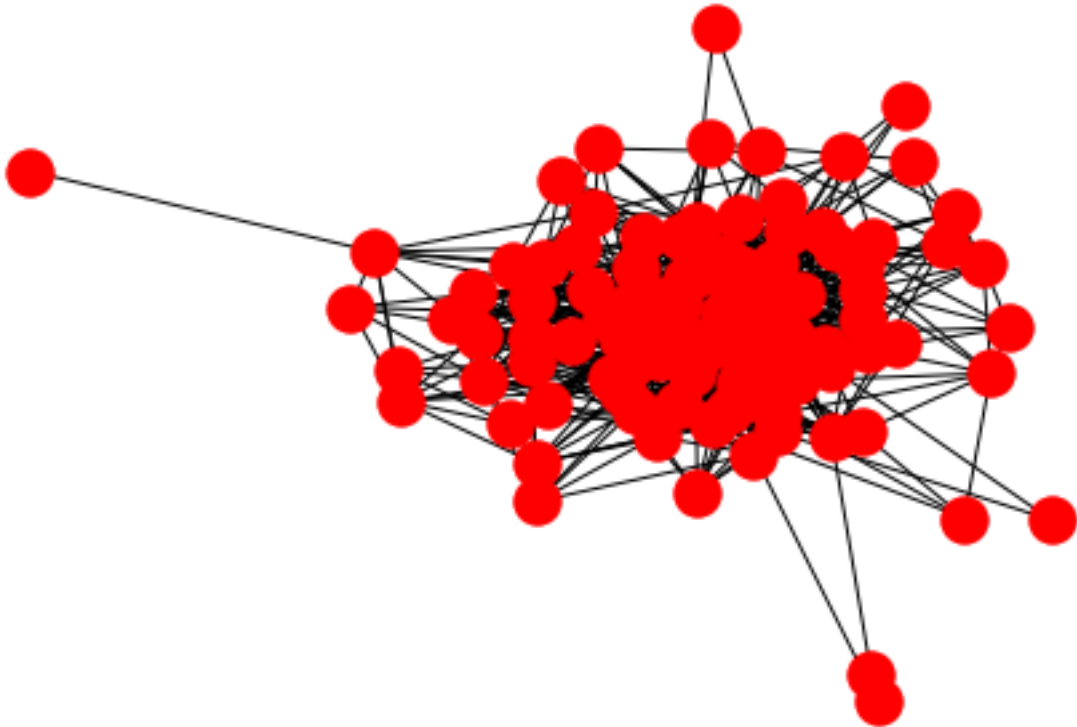
e) One possible way to tap into PageRank is to use the now-familiar package `networkx`. We want you to create a graph from our connections dictionary. Hint: There are many ways of doing this. One possible way involves hooking into the `networkx` method `parse_adjlist` that takes in an adjacency list(connections). Look up documentation on this function to learn more. **(10 points)**

In [336]:

```
{0: [2, 6, 9, 26, 27, 58, 65, 81, 83, 93, 62, 70], 2: [0, 9, 65, 81, 83, 93, 21], 6: [0, 16, 21, 23, 24, 25, 26, 27, 29, 39, 42, 48, 50, 51, 58, 61, 62, 67, 68, 81, 18, 37, 73], 9: [0, 2, 8, 14, 18, 25, 65, 67, 69, 78, 81, 83, 88, 92, 93, 21, 27, 28, 50, 77, 11, 63, 77, 85], 26: [0, 6, 27, 58, 35, 76, 32, 37, 50], 27: [0, 6, 11, 26, 30, 36, 58, 66, 79, 83, 87, 92, 93, 4, 9, 35, 50, 54, 94, 96, 92], 58: [0, 6, 10, 11, 17, 21, 23, 26, 27, 30, 36, 46, 48, 50, 62, 66, 70, 79, 83, 87, 92, 93, 29, 65, 99, 66], 65: [0, 2, 4, 5, 9, 10, 16, 19, 25, 38, 39, 52, 60, 71, 81, 83, 86, 87, 90, 93, 95, 100, 29, 40, 46, 58, 77, 24, 46], 81: [0, 2, 6, 9, 21, 23, 24, 25, 39, 42, 48, 50, 51, 61, 62, 65, 67, 83, 93, 40], 83: [0, 2, 9, 11, 27, 30, 36, 58, 65, 66, 79, 81, 87, 92, 93, 1, 35, 35], 93: [0, 2, 9, 11, 27, 30, 36, 58, 60, 61, 65, 66, 75, 79, 81, 83, 87, 92, 12], 4: [5, 10, 11, 12, 16, 19, 25, 30, 38, 39, 52, 54, 59, 60, 65, 70, 71, 85, 86, 87, 90, 92, 95, 97, 100, 14, 20, 27, 35, 38, 51, 55, 68, 71], 5: [4, 10, 11, 12, 16, 19, 25, 30, 38, 39, 40, 42, 50, 52, 56, 60, 65, 70, 71, 84, 86, 87, 90, 92, 95, 100, 39, 46, 51, 94], 10: [4, 5, 14, 16, 17, 19, 21, 23, 25, 30, 38, 39, 46, 48, 50, 52, 58, 60, 62, 65, 70, 71, 86, 87, 90, 92, 95, 100], 11: [4, 5, 12, 27, 30, 36, 39, 58, 66, 70, 79, 83, 87, 92, 93, 3, 31, 33, 54, 74, 78, 86, 87, 88, 89, 94, 95, 96, 97, 98, 99, 100, 101, 102, 103, 104, 105, 106, 107, 108, 109, 110, 111, 112, 113, 114, 115, 116, 117, 118, 119, 120, 121, 122, 123, 124, 125, 126, 127, 128, 129, 130, 131, 132, 133, 134, 135, 136, 137, 138, 139, 140, 141, 142, 143, 144, 145, 146, 147, 148, 149, 150, 151, 152, 153, 154, 155, 156, 157, 158, 159, 160, 161, 162, 163, 164, 165, 166, 167, 168, 169, 170, 171, 172, 173, 174, 175, 176, 177, 178, 179, 180, 181, 182, 183, 184, 185, 186, 187, 188, 189, 190, 191, 192, 193, 194, 195, 196, 197, 198, 199, 200, 201, 202, 203, 204, 205, 206, 207, 208, 209, 210, 211, 212, 213, 214, 215, 216, 217, 218, 219, 220, 221, 222, 223, 224, 225, 226, 227, 228, 229, 230, 231, 232, 233, 234, 235, 236, 237, 238, 239, 240, 241, 242, 243, 244, 245, 246, 247, 248, 249, 250, 251, 252, 253, 254, 255, 256, 257, 258, 259, 260, 261, 262, 263, 264, 265, 266, 267, 268, 269, 270, 271, 272, 273, 274, 275, 276, 277, 278, 279, 280, 281, 282, 283, 284, 285, 286, 287, 288, 289, 290, 291, 292, 293, 294, 295, 296, 297, 298, 299, 300, 301, 302, 303, 304, 305, 306, 307, 308, 309, 310, 311, 312, 313, 314, 315, 316, 317, 318, 319, 320, 321, 322, 323, 324, 325, 326, 327, 328, 329, 330, 331, 332, 333, 334, 335, 336, 337, 338, 339, 340, 341, 342, 343, 344, 345, 346, 347, 348, 349, 350, 351, 352, 353, 354, 355, 356, 357, 358, 359, 360, 361, 362, 363, 364, 365, 366, 367, 368, 369, 370, 371, 372, 373, 374, 375, 376, 377, 378, 379, 380, 381, 382, 383, 384, 385, 386, 387, 388, 389, 390, 391, 392, 393, 394, 395, 396, 397, 398, 399, 400, 401, 402, 403, 404, 405, 406, 407, 408, 409, 410, 411, 412, 413, 414, 415, 416, 417, 418, 419, 420, 421, 422, 423, 424, 425, 426, 427, 428, 429, 430, 431, 432, 433, 434, 435, 436, 437, 438, 439, 440, 441, 442, 443, 444, 445, 446, 447, 448, 449, 450, 451, 452, 453, 454, 455, 456, 457, 458, 459, 460, 461, 462, 463, 464, 465, 466, 467, 468, 469, 470, 471, 472, 473, 474, 475, 476, 477, 478, 479, 480, 481, 482, 483, 484, 485, 486, 487, 488, 489, 490, 491, 492, 493, 494, 495, 496, 497, 498, 499, 500, 501, 502, 503, 504, 505, 506, 507, 508, 509, 510, 511, 512, 513, 514, 515, 516, 517, 518, 519, 520, 521, 522, 523, 524, 525, 526, 527, 528, 529, 530, 531, 532, 533, 534, 535, 536, 537, 538, 539, 540, 541, 542, 543, 544, 545, 546, 547, 548, 549, 550, 551, 552, 553, 554, 555, 556, 557, 558, 559, 560, 561, 562, 563, 564, 565, 566, 567, 568, 569, 570, 571, 572, 573, 574, 575, 576, 577, 578, 579, 580, 581, 582, 583, 584, 585, 586, 587, 588, 589, 590, 591, 592, 593, 594, 595, 596, 597, 598, 599, 600, 601, 602, 603, 604, 605, 606, 607, 608, 609, 610, 611, 612, 613, 614, 615, 616, 617, 618, 619, 620, 621, 622, 623, 624, 625, 626, 627, 628, 629, 630, 631, 632, 633, 634, 635, 636, 637, 638, 639, 640, 641, 642, 643, 644, 645, 646, 647, 648, 649, 650, 651, 652, 653, 654, 655, 656, 657, 658, 659, 660, 661, 662, 663, 664, 665, 666, 667, 668, 669, 670, 671, 672, 673, 674, 675, 676, 677, 678, 679, 680, 681, 682, 683, 684, 685, 686, 687, 688, 689, 690, 691, 692, 693, 694, 695, 696, 697, 698, 699, 700, 701, 702, 703, 704, 705, 706, 707, 708, 709, 710, 711, 712, 713, 714, 715, 716, 717, 718, 719, 720, 721, 722, 723, 724, 725, 726, 727, 728, 729, 730, 731, 732, 733, 734, 735, 736, 737, 738, 739, 740, 741, 742, 743, 744, 745, 746, 747, 748, 749, 750, 751, 752, 753, 754, 755, 756, 757, 758, 759, 760, 761, 762, 763, 764, 765, 766, 767, 768, 769, 770, 771, 772, 773, 774, 775, 776, 777, 778, 779, 780, 781, 782, 783, 784, 785, 786, 787, 788, 789, 790, 791, 792, 793, 794, 795, 796, 797, 798, 799, 800, 801, 802, 803, 804, 805, 806, 807, 808, 809, 810, 811, 812, 813, 814, 815, 816, 817, 818, 819, 820, 821, 822, 823, 824, 825, 826, 827, 828, 829, 830, 831, 832, 833, 834, 835, 836, 837, 838, 839, 840, 841, 842, 843, 844, 845, 846, 847, 848, 849, 850, 851, 852, 853, 854, 855, 856, 857, 858, 859, 860, 861, 862, 863, 864, 865, 866, 867, 868, 869, 870, 871, 872, 873, 874, 875, 876, 877, 878, 879, 880, 881, 882, 883, 884, 885, 886, 887, 888, 889, 890, 891, 892, 893, 894, 895, 896, 897, 898, 899, 900, 901, 902, 903, 904, 905, 906, 907, 908, 909, 910, 911, 912, 913, 914, 915, 916, 917, 918, 919, 920, 921, 922, 923, 924, 925, 926, 927, 928, 929, 930, 931, 932, 933, 934, 935, 936, 937, 938, 939, 940, 941, 942, 943, 944, 945, 946, 947, 948, 949, 950, 951, 952, 953, 954, 955, 956, 957, 958, 959, 960, 961, 962, 963, 964, 965, 966, 967, 968, 969, 970, 971, 972, 973, 974, 975, 976, 977, 978, 979, 980, 981, 982, 983, 984, 985, 986, 987, 988, 989, 990, 991, 992, 993, 994, 995, 996, 997, 998, 999]}
```

In [302]:

In [303]:



f) Use the template from the midterm and previous assignments, or just use networkx's `pagerank()` API that professor mentioned in class to calculate pageRank for google employees. Using this information, find the 10 most central googlers, and the 10 least central googlers. Are they the same as the 10 most and 10 least connected googlers? **(10 points)**

In [304]:

In [305]:

In [306]:

Out[306]:

```
[ ('50', 0.022820129402410672),  
  ('60', 0.02194001836738658),  
  ('52', 0.021460032170229196),  
  ('67', 0.020625117069219122),  
  ('25', 0.02033308102254799),  
  ('16', 0.020250826019185545),  
  ('39', 0.018278153501639652),  
  ('30', 0.018227946149985987),  
  ('71', 0.016717148421104666),  
  ('4', 0.01644505221227893) ]
```

In [307]:

Out[307]:

```
[(50, 52),  
 (16, 46),  
 (52, 45),  
 (25, 43),  
 (39, 43),  
 (60, 42),  
 (67, 41),  
 (30, 39),  
 (71, 38),  
 (48, 36)]
```

In [308]:

Out[308]:

```
[('99', 0.004468905019286139),  
 ('91', 0.004172302704734399),  
 ('94', 0.003822089001090339),  
 ('64', 0.0035784682084692614),  
 ('57', 0.002957342386297933),  
 ('80', 0.002620839816228804),  
 ('15', 0.002409624501967589),  
 ('53', 0.002226273708244273),  
 ('44', 0.002080589913150963),  
 ('72', 0.001996646491383476)]
```

In [309]:

Out[309]:

```
[(97, 6),
 (91, 5),
 (94, 5),
 (64, 4),
 (57, 3),
 (15, 2),
 (80, 2),
 (44, 1),
 (53, 1),
 (72, 1)]
```

Note: We can see that the top and bottom pagerank and connected employee list are different (as we have noticed in class before) in that the top members of pagerank are not the same as the top members with the most connections. Those with connections indirectly others with MORE connections often increases the pagerank more - although there is much overlap between the top ten groups. Given the analogy of Sergey Brin previously - a connection to Sergey Brin at this point raises a person's pagerank more than a connection to a co-op student (assuming the co-op student has very few connections in turn) so it is more of a measure of importance!

g) BONUS: Can you plot google's connectome? That connectome should easily reveal the least central googlers, since they represent isolated nodes, far away from network traffic. You have some latitude on how you choose to do this. Can you do something to upend that graph to also easily reveal the most central googlers? **(5 points)**

In []:

Solutions: Part 2

a) Write a list comprehension (ideally) or some type of loop to determine connections of connections. Essentially, you should already have a list of connections, but now search the connections of a specific googler to find connections of connections. For example, if (Ian Isla, Kaleb Lauren) were connected in Part 1, and (Kalen Lauren, Roman Valentina) were also connected, then Roman Valentina would be a "connection of a connection" of Ian Isla. While you do this, perform a check to ensure that the connection of a connection should not already be a connection itself. **(20 points)**

In [310]:

Out[310]:

	0	1	2	3	4	5	6	7	8	9	...	42	43	44	45	46	47
0	2	6.0	9.0	26.0	27.0	58.0	65.0	81.0	83.0	93.0	...	NaN	NaN	NaN	NaN	NaN	NaN
2	0	9.0	65.0	81.0	83.0	93.0	21.0	NaN	NaN	NaN	...	NaN	NaN	NaN	NaN	NaN	NaN
6	0	16.0	21.0	23.0	24.0	25.0	26.0	27.0	29.0	39.0	...	NaN	NaN	NaN	NaN	NaN	NaN
9	0	2.0	8.0	14.0	18.0	25.0	65.0	67.0	69.0	78.0	...	NaN	NaN	NaN	NaN	NaN	NaN
26	0	6.0	27.0	58.0	35.0	76.0	32.0	37.0	50.0	NaN	...	NaN	NaN	NaN	NaN	NaN	NaN

5 rows × 52 columns

In [311]:

In [312]:

Out[312]:

	0	1	2	3	4	5	6	7	8	9	...	43	44	45	46	47	48
0	2	6.0	9.0	26.0	27.0	58.0	65.0	81.0	83.0	93.0	...	NaN	NaN	NaN	NaN	NaN	NaN
1	0	9.0	65.0	81.0	83.0	93.0	21.0	NaN	NaN	NaN	...	NaN	NaN	NaN	NaN	NaN	NaN
2	0	16.0	21.0	23.0	24.0	25.0	26.0	27.0	29.0	39.0	...	NaN	NaN	NaN	NaN	NaN	NaN
3	0	2.0	8.0	14.0	18.0	25.0	65.0	67.0	69.0	78.0	...	NaN	NaN	NaN	NaN	NaN	NaN
4	0	6.0	27.0	58.0	35.0	76.0	32.0	37.0	50.0	NaN	...	NaN	NaN	NaN	NaN	NaN	NaN

5 rows × 53 columns

In [339]:

Count of Secondary Connections: 3214

In [334]:

```
{0: [4, 5, 10, 11, 12, 16, 19, 25, 30, 38, 39, 52, 54, 60, 71, 85, 86,
87, 90, 92, 95, 100, 40, 42, 50, 84, 21, 23, 24, 29, 48, 51, 61, 67, 6
8, 7, 8, 14, 28, 69, 18, 78, 88, 17, 46, 36, 66, 79, 75, 1, 35, 37, 73
, 77, 96, 22, 64, 99, 32, 63, 76, 94], 2: [6, 26, 27, 58, 4, 5, 10, 11
, 12, 16, 19, 25, 30, 38, 39, 52, 60, 70, 71, 85, 86, 87, 90, 92, 95,
97, 100, 40, 42, 50, 23, 24, 29, 48, 51, 61, 62, 67, 8, 14, 28, 69, 18
, 78, 88, 17, 46, 36, 66, 79, 75, 1, 35, 77, 72, 63], 6: [9, 65, 83, 9
3, 4, 5, 10, 11, 12, 19, 30, 38, 52, 54, 60, 70, 71, 85, 86, 87, 90, 9
2, 95, 97, 100, 40, 56, 84, 7, 8, 14, 20, 28, 31, 55, 69, 78, 88, 17,
46, 36, 66, 79, 13, 41, 34, 75, 1, 3, 74, 35, 77, 96, 72, 22, 43, 45,
64, 99, 32, 63, 76, 94, 89, 47, 49, 98, 57], 9: [26, 58, 4, 5, 10, 12,
16, 19, 30, 38, 39, 52, 54, 59, 60, 70, 71, 86, 87, 90, 95, 97, 100, 4
0, 42, 56, 84, 23, 24, 29, 48, 51, 61, 62, 68, 7, 20, 31, 55, 17, 46,
36, 66, 79, 13, 34, 75, 1, 3, 33, 74, 35, 37, 73, 96, 72, 22, 43, 45,
64, 99, 76, 94, 89, 47, 98, 44], 26: [65, 81, 83, 93, 4, 5, 10, 11, 12
, 16, 25, 30, 39, 52, 54, 60, 70, 71, 85, 87, 92, 40, 42, 56, 84, 21,
23, 24, 29, 48, 51, 61, 62, 67, 68, 7, 8, 14, 20, 31, 55, 69, 18, 78,
17, 46, 36, 66, 79, 75, 1, 74, 73, 96, 99, 63, 94, 89, 98, 53, 91], 27
: [65, 81, 5, 10, 12, 16, 19, 25, 38, 39, 52, 59, 60, 70, 71, 85, 86,
```

b) What's the average number of connections-of-connections at google? You will be surprised to find out that even with only 3,000 employees, and just a few projects, cruises, and sports in common, the number of connections² is huge. So, it would take forever to iterate through all possible employees. DON'T DO THIS. Instead, just use a random sampling technique to *sample* through a fraction of employees, throw in some debugging print statements in your loop(s) to see how fast your sampling is progressing, produce a few random samplings and average over these samplings. We leave it up to you how many samples you decide to choose, and will be lenient with our grading as long as your procedure seems reasonable. **(10 points)**

In [340]:

Average Secondary Connections (Amongst Everyone): 32.14

Average Secondary Connections (Only Amongst Googlers With Secondary Connections): 32.14

Note: in this case everyone has a secondary connection - but that is not necessarily true!

c) BONUS: What's the average distance that separates any two googlers? You may want to read [this](http://blogs.cornell.edu/info2040/2011/09/26/networks-in-hollywood-the-six-degrees-of-kevin-bacon/) (<http://blogs.cornell.edu/info2040/2011/09/26/networks-in-hollywood-the-six-degrees-of-kevin-bacon/>) article about Kevin Bacon. [This](https://en.wikipedia.org/wiki/Erd%C5%91s%E2%80%93Bacon_number) (https://en.wikipedia.org/wiki/Erd%C5%91s%E2%80%93Bacon_number) is a funny article, too. **(5 points)**

In []: