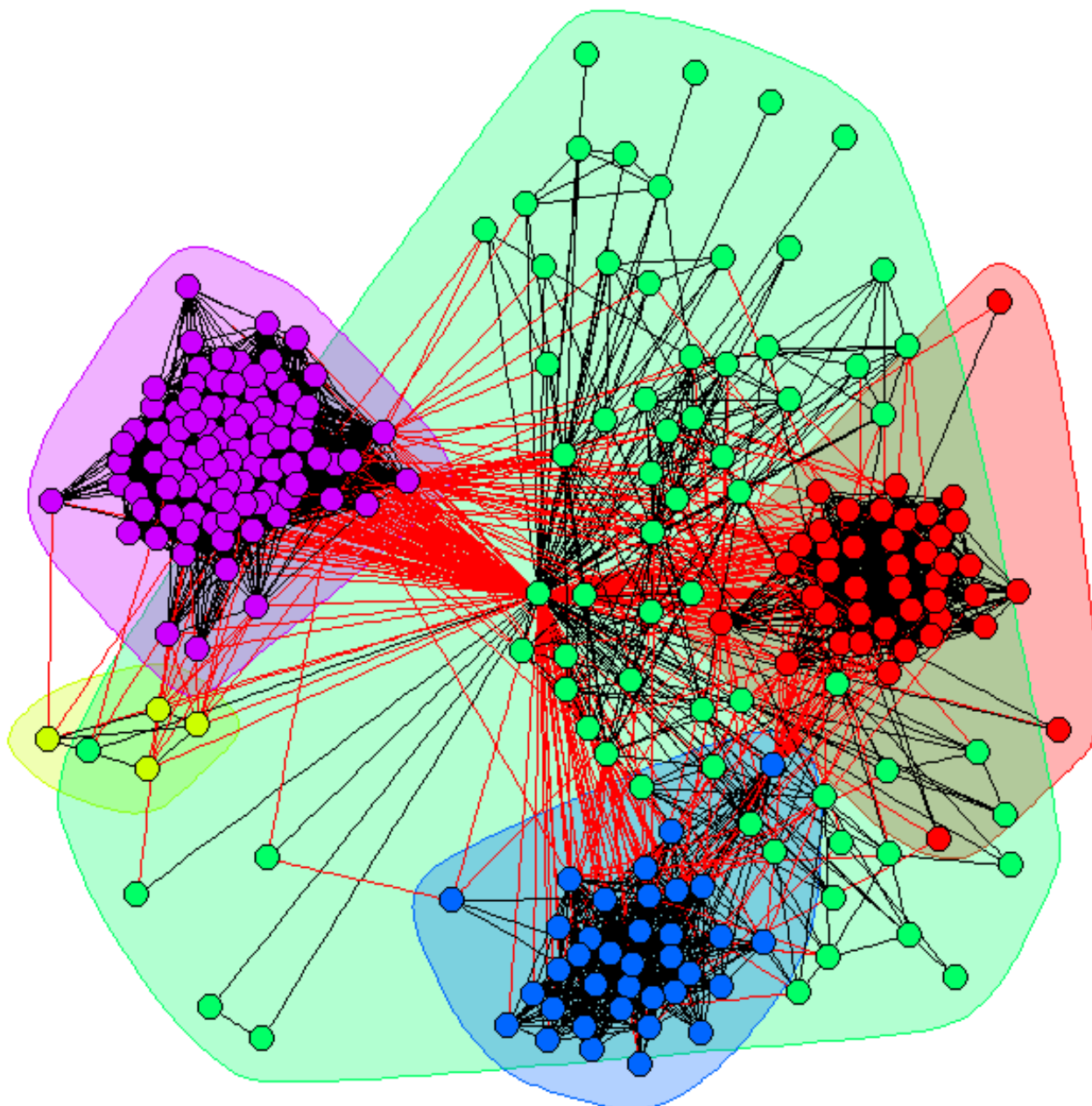# Community Detection in Social Networks

🕑 October 19, 2014     📁 generic ML     🏷 Python, R, social network analysis

In this post I would like to share a very basic approach to **Community Detection in Social Networks**.

I came across this fascinating topic following the superb course on **Mining Massive Datasets** provided on Coursera by Stanford University. The specific field of finding overlapping clusters in graphs is introduced and deeply treated during the third week of classes (links to the PDF slides available for Part 1 and Part 2). I immediately found it extremely interesting and decided to play around by myself. There are at least two very strong reasons to directly check the potentials of these group of algorithms: first of all my complete lack of knowledge in the field and secondly the data I found a couple of weeks ago on a the Kaggle competition "Learning Social Circles in Networks". The contest challenges participants to correctly infer Facebook users' social communities. Such circles may be disjoint, overlap, or be hierarchically nested.  To do this, machine learners have access to:

1. A list of the user's friends
2. Anonymized Facebook profiles of each of those friends
3. A network of connections between those friends (their "ego network")

This is exactly what I needed for my learning purposes!

# Overview

The approach I propose below is structured in two main parts:

1. **Build the Graph** of the ego-networks extracting nodes and edges from Kaggle data. I implemented this step in Python, generating the graphs with Networkx and saving the Adjiacency matrix of each of them to a separate file.
2. **Community Detection** on top of the undirected graph. I performed this step in R, loading the graphs as Adjiacency matrices and then run a bunch of Clustering Algorithms available in R-igraph.

The use of both Python and R was not planned in the first place. I directly dived into the first of them supported by Neyworkx, but as soon as I started deepening the community detection algorithms I realized that R-igraph had a woderful ensemble of methods directly available. Note that igraph supports Python as well but apparently there are not the same features between the two libraries and the R one seems to be much fancier. I was a bit disappointed at the very beginning but in the end I grabbed the opportunity of learning a new package.

Enough words, I'd say. Let's go for some code.

# Building Ego-Networks

The Kaggle data (available here) is organized in 110 .egonet files (corresponding to 110 anonymized facebook users), each containing the network of his friends. A practicle example may help to clarify the data structure.

Let's focus on the file **0.egonet**, which contains all the information on **user 0**'s network. Each row of the file is the list of the friends of the first user in the line who is directly part of the ego's network. Below the first 4 lines are shown (for the purpose of clearness only the first 5 five connections in the line are reported).

```
1  1: 146 189 229 201 204 ...
2  2: 146 191 229 201 204 ...
3  3: 185 80 61 188 22 222 ...
4  4: 72 61 187 163 177 138 ...
```

**0** has **1** as friend who has **146-189-229**… as friends as well.

**0** has **2** as friend who has **146-191-229**… as friends as well.

**0** has **3** as friend who has **185-80-61**… as friends as well.

**0** has **4** as friend who has **72-61-187**… as friends as well.

Well I guess you get the point…

Below  I attach the **Python code** which access every egonet file and builds a list of nodes and edges to be fed to the Networkx constructor.  Just to be clear [0, 1, 2, 3, 4 …] are vertices of the graph while [(0-1), (1-146), (1-189), (1-229) …] are edges or connections. After a graph has been constructed its adjiacency

matrix is computed and saved in a csv file.

```python
1  import networkx as nx
2  from os import listdir
3  from os.path import isfile, join
4  import itertools
5  import matplotlib.pyplot as plt
6  import os
7  import re
8  import scipy
9  from scipy.sparse import *
10 from operator import itemgetter
11 from sklearn.cluster import KMeans
12 import numpy as np
13 import sys
14 import pandas as pd
15
16
17 def load_egonet_files(path):
18     """
19     given the path to the .egonet files returns a list with all the files.
20     """
21     onlyfiles = [fyle for fyle in listdir(path) if fyle.endswith('.egonet')]
22     return onlyfiles
23
24 ###############################################################################
25
26 def build_graph(n):
27     """
28     In the Kaggle Competition Learning Social Circles in Networks (https://www.kag
29     110 ego-ketworks are available. The function takes as argument an integer in t
30     a networkx graph of the friends network of the user.
31     """
32     edges = []
33     nodes = []
34     path = 'D:\KaggleGraphs\egonets\egonets'
35     egonets = load_egonet_files(path)
36     for egonet in egonets[n:n+1]:
37         ego = int(re.match( r'([0-9]+).egonet', egonet).group(1))
38         m = open(os.path.join(path,egonet), "r")
39         friends = [line[:-1].replace(':','').split(' ') for line in m.readlines()]
40         friends = [map(int, friend[:1]) if friend[1] == '' else map(int, friend) f
41         edges += [(ego,friend[0]) for friend in friends]
42
43         for friend in friends:
44             edges += [(friend[0], user) for user in friend[1:] if len(friend)>1]
45
46         nodes += list(itertools.chain.from_iterable(friends)) + [ego]
47     edges = list(set(tuple(sorted(edge)) for edge in edges))
48     nodes = list(set(nodes))
49     G = nx.Graph()
50     G.add_nodes_from(nodes)
51     G.add_edges_from(edges)
52
53     return G, nodes
54
55 ###############################################################################
56
57 def to_R(n):
58     """
59     generates a CSV file with the Adjacency matrix representation of G.
60     """
61     G, nodes = build_graph(n)
62     A = nx.to_numpy_matrix(G)
```

```
63        df = pd.DataFrame(A, index=nodes, columns=nodes)
64        df.to_csv('graph-{}.csv'.format(n),index=True)
65
66  ###############################################################################
67
68  if __name__ == '__main__':
69
70      for n in range(110):
71          to_R(n)
```

The result of the provided code are 110 CSV files containing the adjiacency matrices of each ego network graph. Let's move to the real Clustering part.
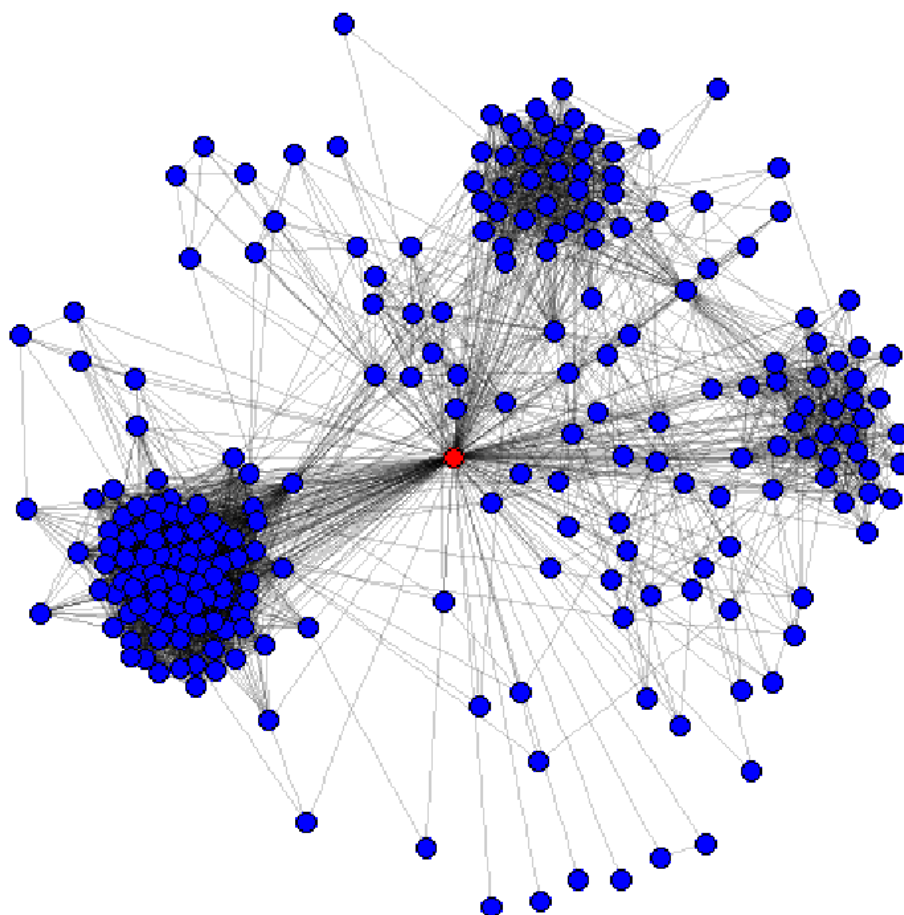
# Detecting Communities

First of all let's plot a graph and see how it looks like before clustering detection. Below the **R code** to load the data from CSV file, build the network (we stick to the 0.egonet) and draw it.

```
1  library(igraph)
2  # read graph from csv file
3  dat = read.csv('graph-0.csv', header=TRUE, row.names=1, check.names=FALSE)
4  m = as.matrix(dat)
5  # build graph from adjacency matrix
6  g = graph.adjacency(m,mode="undirected",weighted=NULL)
7
8  # plots the graph
9  E(g)$color <- rgb(0,0,0,alpha=.2)
10 ego <- names(which.max(degree(g)))
11 V(g)[V(g) != ego]$color = 'blue'
12 V(g)[ego]$color = 'red'
13
14 windows()
15 plot(g, vertex.label=NA, vertex.size=5, layout=layout.fruchterman.reingold)
```

Time for some clustering.

R-igraph provides several powerful community detection algorithms. Each of them works in a different way and I highly encourage you to have a look at this very informative post on Stack Overflow describing all of them in detail. I decided to go for the whole bunch of algorithms as I wanted to somewhat compare their performances, which I did with the help of **Modularity**. This metrics measures the strength of division of a network into modules. Networks with high modularity have dense connections between the nodes within modules but sparse connections between nodes in different modules.

Modularity is basically the fraction of the edges that fall within the given groups minus the expected such fraction if edges were distributed at random. So the higher the better.

Here you find the results on the user-0-network.

```
1  > wc <- walktrap.community(g)
2  > modularity(wc)
```
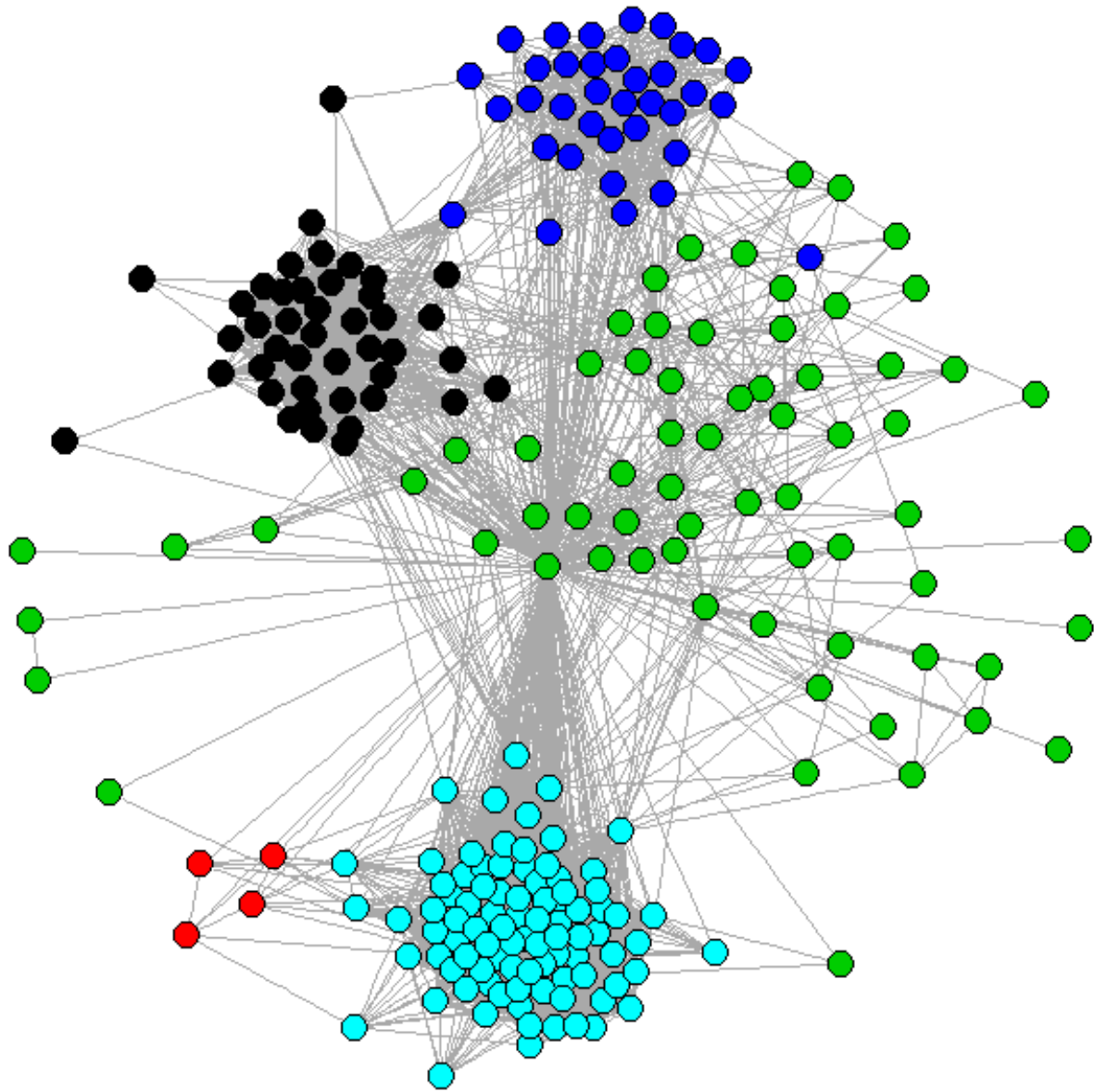
```
3   [1] 0.4629543
4
5   > wc <- fastgreedy.community(g)
6   > modularity(wc)
7   [1] 0.4463902
8
9   > wc <- edge.betweenness.community(g)
10  > modularity(wc)
11  [1] 0.4330911
12
13  > wc <- spinglass.community(g)
14  > modularity(wc)
15  [1] 0.4649535
16
17  > wc <- leading.eigenvector.community(g)
18  > modularity(wc)
19  [1] 0.4511259
20
21  > wc <- label.propagation.community(g)
22  > modularity(wc)
23  [1] 0.4314803
```

The spinglass.community algorithm (based on a statistical physics approach) is the best one, with a modularity of 0.4649. Turns out that for this particular problem of community detection in small ego-social-networks the spinglass method beats the others in all the 110 egonet graphs.

Below you can find a nice visualization of the detected clusters, in R as well. By the way the plot at the top of the post is exactly the same as the following one visualized in a fancier way.

```
1   plot(g, vertex.label=NA, vertex.size=5, vertex.color=membership(wc),
2       layout=layout.fruchterman.reingold)
```

Share! 🅕 ⟨0⟩  🅖+ ⟨0⟩  🅛🅽 ⟨3⟩  🅥 ⟨1⟩  ✉

by *Francesco Pochetti*

# Comments

0 comments