

ADM – Homework4

Giacomo Parmendola(1462237)

Vigèr Durand Azimedem Tsafack(1792126)

Mattia Podio(1554740)

✓ First Step: Graph creation

1. Here we process the JSON file and create a graph with all the authors as nodes. Two nodes are connected if they share, at least, one publication and the weight of an edge is computed in this way: $w(a1,a2) = 1 - J(p1,p2)$.

Find the python code in the file *graph_creation.py*.

See the explanation of the procedure used in the (.rm) file on the GitHub repository.

Output:

Name:

Type: Graph

Number of nodes: 904664

Number of edges: 3679276

Average degree: 8.1340

Computation time:

...data loaded...

...graph creation completed...

Elapsed time: 95.88851046562195

...weights computation completed...

Elapsed time: 42.06702256202698

✓ Second Step: Some statistics

1. Given a conference in input, return the subgraph induced by the set of authors who published at the input conference at least once: Find the python code in the file

graph_creation.py

First, we ask the user to insert some conference ID:

```
*****
Given a conference in input, the program returns the subgraph induced
by the set of authors who published at the input conference at least once
*****
```

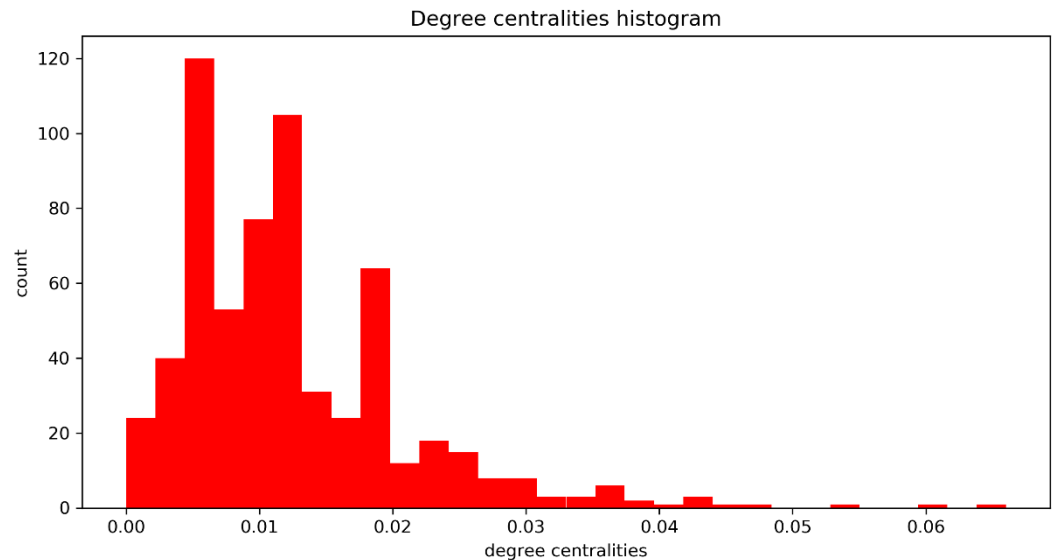
Insert the conference ID: 3345

Computation time:

...first subgraph creation completed...

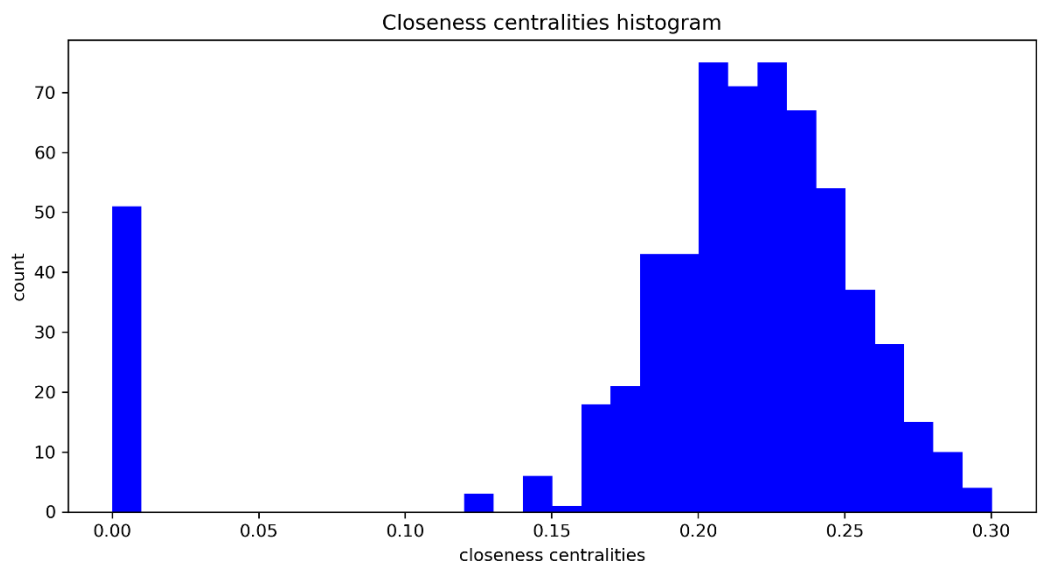
Elapsed time: 0.280933141708374

 Here are some statistics



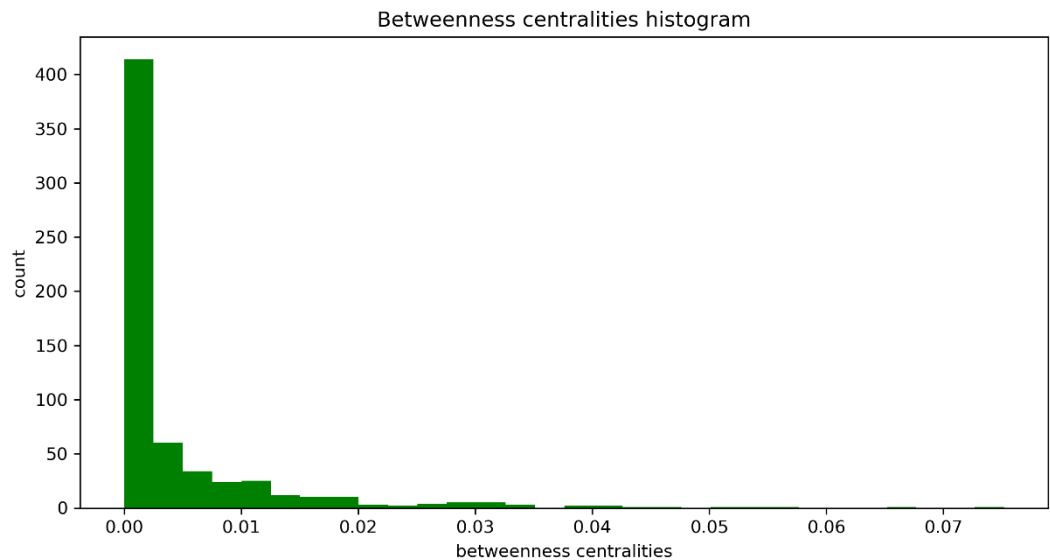
Degree is a simple centrality measure that counts how many neighbors a node has. Thus, A node is important if it has many neighbors. From this histogram, we can state that the probability for a node to be important is quite low while the probability for a node to be less important is pretty high. (considering count as probability measure)

degree centrality = 0.005 → count = 120 → low importance, high probability
 degree centrality = 0.065 → count = 1 → high importance, low probability



Closeness centrality measures the mean distance from a vertex to other vertices. it gives low values to more central nodes and high values to less central ones, which is the opposite of other centrality measures. Here the histogram shows that the probability for a node to be important according to the closeness centrality is low while the probability to be less important is high.

closeness centrality = 0.001 → count = 51 → high importance, medium probability
 closeness centrality = 0.24 → count = 72 → low importance, high probability



Betweenness centrality measures the extent to which a vertex lies on paths between other vertices. Vertices with high betweenness may have considerable influence within a network by virtue of their control over information passing between others. They are also the ones whose removal from the network will most disrupt communications between other vertices. Here looking at the histogram we can say that we are again in a case where the probability for a node to be important is quite low while the probability to be less important is pretty high.

betweenness centrality = 0.001 → count = 425 → low importance, high probability
 betweenness centrality = 0.04 → count = 1 → high importance, low probability

We can generally state that two over three centrality measures (Degree and betweenness) follow the power law (Pareto principle also known as the 80/20 rule or the law of the vital few) which state that for many events, roughly 80% of the effects come from 20% of the causes. For example, 80% of Italy's land is owned by 20% of the population.

Computation time:

```
...centrality measures computation done...
...histograms creation done...
Elapsed time: 2.583315849304199
```

- given in input an author and an integer d , get the subgraph induced by the nodes that have hop distance at most equal to d with the input author:
 find the python code in the file [graph_creation.py](#).
 find the `bfsr()` and `bfsr()` functions used here in the file [Libhw4.py](#)

First, we ask the user to insert some author ID:

```
*****
Given in input an author and an integer d the program returns the subgraph
induced by the nodes that have hop distance at most equal to d
*****
```

Insert an author ID: 16749

We also ask the user for an integer d which is going to be the max for the hop-distance:

Insert an integer d: 2

Then, we ask again the user to tell whether he wishes to use the recursive, the iterative breath first search algorithm (see the explanation of both algorithms in the (.rm) file on the GitHub repository) or even the *networkx.ego_graph()* to compute the hop-distances:

```
Choose what function you wish to use in order to compare the computation times
1 for iterative breath first search algorithm
2 for recursive breath first search algorithm
3 for networkx.ego_graph()
  essentially based on single source Dijkstra
e to exit
```

1

```
...subgraph creation completed...
Elapsed time: 0.0
```

2

```
...subgraph creation completed...
Elapsed time: 0.0009760856628417969
```

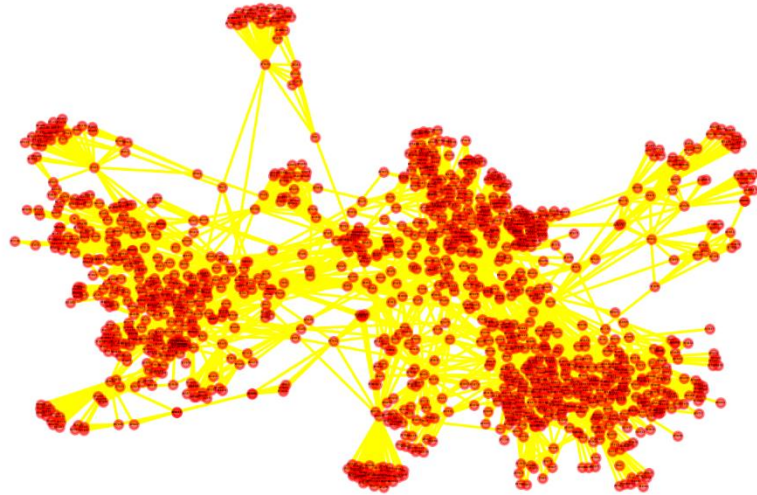
3

```
...subgraph creation completed...
Elapsed time: 0.016414880752563477
```

The outputs show that all the three algorithms have similar behavior in computation time. Thus, the reason why we decided to implement an iterative version of the *bfs* was indeed to avoid raising the *RecursionError* while going deeper into the graph.

Output graph:

```
*****  
Now let's visualise the subgraph  
*****
```



✓ **Third Step:** Erdős number

1. Here we compute the weight of the shortest path that connects some input author with Aris. We use as measure of distance the weight of the edges.
Find the python code in the file *graph_creation.py*
Find the *shortest_path()* and *our_dijkstra()* functions used here in the file *Libhw4.py*

First, we ask the user to insert some author ID:

Insert the author ID: 16404

We then ask again the author to choose which function he wishes to use between *shortest_path()* (adaptation of the bfs algorithm transforming it into a sort of dijkstra algorithm) and *nx.dijkstra_path_length()*:

```
Choose what function you wish to use in order to compare the computation times  
1 for Libhw4.shortest_path()  
2 for nx.dijkstra_path_length()  
e to exit
```

1

Shortest path between Aris and 16404 is 3.743371099185053

```
...shortest path weight calculation completed...  
Elapsed time: 13.403905153274536
```

2

Shortest path between Aris and 16404 is 3.743371099185053

```
...shortest path weight calculation using the networkx function completed...  
Elapsed time: 11.474873542785645
```

Here we notice that the two computation times are pretty close. If we neglect the availability of the processor while running these two functions, we could say that our algorithm is almost as fast as the networkx implementation of the Dijkstra algorithm.

```
In [78]: timeit nx.dijkstra_path_length(G, 16404, 256176)  
16.9 s ± 3.64 s per loop (mean ± std. dev. of 7 runs, 1 loop each)
```

```
In [79]: timeit lb.shortest_path(G, 16404, 256176)  
12.9 s ± 269 ms per loop (mean ± std. dev. of 7 runs, 1 loop each)
```

2. Here we compute for each node of the graph, its GroupNumber, defined as follow:

$GroupNumber(v) = \min_{u \in \{ShortestPath(v,u)\}}$.

Find the python code in the file *graph_creation.py*

Find the function *GroupNumbers()* in the file *Libhw4.py*

First, we ask the user to insert a list of author ID's:

```
Insert a list of author ID's separated by spaces  
max 21 items: 365066 273515 17528 19210 364934
```

We then call the *GroupNumbers()* function which returns:

```

904640: (3.8416414141414146, 364934),
904641: (4.399658149204429, 365066),
904642: (3.880730427764326, 17528),
904643: (3.9526760590036605, 19210),
904644: (3.9526760590036605, 19210),
904645: (3.9526760590036605, 19210),
904646: (3.9526760590036605, 19210),
904647: (4.546729362591432, 17528),
904648: (4.380696261293276, 17528),
904649: (4.380696261293276, 17528),
904650: (4.380696261293276, 17528),
904651: (3.8506427515787736, 364934),
904652: (4.721212302466401, 273515),
904653: (4.918326246340947, 19210),
904654: (4.918326246340947, 19210),
904655: (4.918326246340947, 19210),
904656: (4.918326246340947, 19210),
904657: (5.382803209620128, 17528),
904658: (5.382803209620128, 17528),
904659: (4.651683663916531, 273515),
904660: (4.4212210468982915, 365066),
904661: (inf, None),
904662: (inf, None),
904663: (inf, None),
904664: (inf, None)}
Author: (Min Shortest Path, Best Input Node)

```

(inf, None) represent the case when none of the input nodes is reachable from the corresponding node: the min shortest path is *infinity* and the best input node is *None*.

Computation time:

```

...GroupNumber's computation completed...
Elapsed time: 241.12900519371033

```