

Introduction to Network Analysis

Homework 2

Nace Gorenc

April 2021

1 Where is SN100?

The importance of the dolphin SN100 is in this case measured by betweenness centrality since it highlights nodes that are bridges between other nodes. Table 1 represents 5 most important nodes in this network. As it can be observed, there are 2 dolphins with the importance higher than 0.2. Note that the betweenness centrality measure is calculated by the algorithm in `networkx` library, which has slightly different normalization factor ($2/((n-1)(n-2))$) compared to the normalization factor which we have seen in the lectures ($1/n^2$).

The printout of the code used for this problem can be seen below.

```
1 def e1():
2     dolphins = nx.read_adjlist('data/dolphins')
3     bc = nx.betweenness centrality(dolphins, normalized=True)
4     print('Betweenness centrality:')
5     print({key: value for key, value in sorted(bc.items(),
        ↪ key=lambda item: item[1], reverse=True)[:5])})
```

Dolphin's name	Betweenness centrality measure
SN100	0.248
Beescratch	0.213
SN9	0.143
SN4	0.139
DN63	0.118

Table 1: Most important dolphins according to betweenness centrality measure.

2 Is software scale-free?

Figure 1 and 2 represent degree distributions p_k vs. node degree k graphs. In both network the the degree distributions and out-degree distributions have

similar shape until $k < 20$ and after that we can see the similarity between degree and in-degree distributions. The major difference between this two network is that Java language has fewer root libraries that do not depend on any other libraries, however there are a lot of libraries that depends on a small number of other so called root libraries, and Lucene search engine has a lot of root libraries and libraries that depend on fewer than 10 other libraries. Additionally, both networks contain a great number of higher level libraries that no other library depend on.

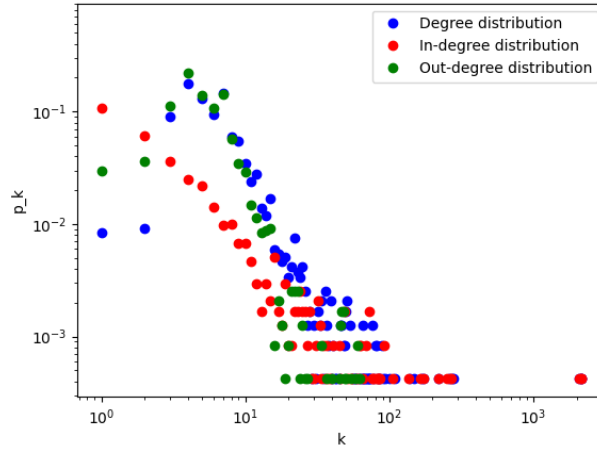


Figure 1: Degree distributions for Java language.

It appears that in-degree distributions of both networks follow a power-law due to linearity of the dots on the logarithmic graph. On the other hand, degree and out-degree distributions of both graph seem to follow Poisson distribution since the shape of the graphs resemble the Poisson distribution.

Using the likelihood formula, we can plot a graph γ vs. K_{min} of both in-degree distributions. As we can observe, both plots in figures 3 and 4 flatten when k_{min} is greater than 25 thus this is our k_{min} . Furthermore, the power-law exponents γ of Java and Lucene network for $k_{min} = 25$ are 2.14 and 2.12, respectively.

The printout of the code used for this problem can be seen below.

```
1 def generate_plot(s):
2     g = read_undirected('data/' + s)
3     tmp = [x for _, x in g.degree()]
4     p_k = [tmp.count(i) / g.number_of_nodes() for i in
5            range(max(tmp) + 1)]
6     plt.plot(p_k, 'bo')
7     tmp = [x for _, x in g.in_degree()]
```

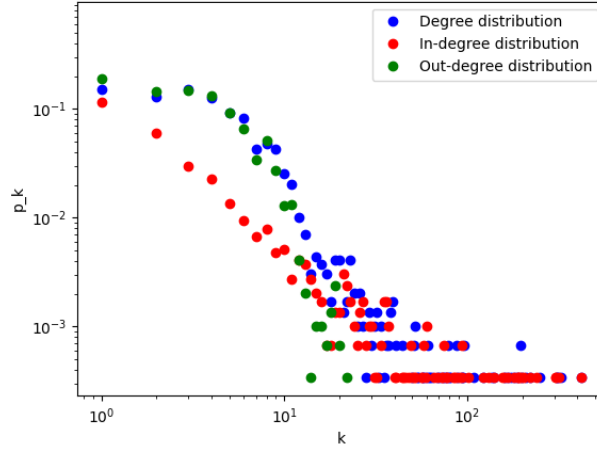


Figure 2: Degree distributions for Lucene search engine.

```

7     p_k = [tmp.count(i) / g.number_of_nodes() for i in
    ↪     range(max(tmp) + 1)]
8     plt.plot(p_k, 'ro')
9     tmp = [x for _, x in g.out_degree()]
10    p_k = [tmp.count(i) / g.number_of_nodes() for i in
    ↪     range(max(tmp) + 1)]
11    plt.plot(p_k, 'go')
12    plt.yscale('log')
13    plt.xscale('log')
14    plt.legend(['Degree distribution', 'In-degree distribution',
    ↪     'Out-degree distribution'],
15              loc='upper right')
16    plt.xlabel('k')
17    plt.ylabel('p_k')
18    plt.savefig(s + '.png')
19    plt.show()
20
21
22    def gamma(s):
23        k = [i for i in range(1, 101)]
24        gamma = []
25        g = read_undirected('data/' + s)
26        tmp = [x for _, x in g.in_degree()]
27        for k_min in k:
28            n = sum([i >= k_min for i in tmp])
29            gamma.append(1 + n / (sum([np.log(i / (k_min - 0.5)) if i >=
    ↪            k_min else 0 for i in tmp])))

```

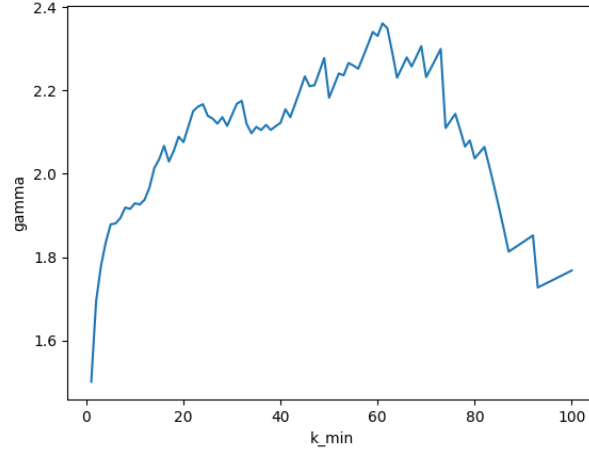


Figure 3: Java language.

```

30 plt.plot(k, gamma)
31 plt.xlabel('k_min')
32 plt.ylabel('gamma')
33 plt.savefig(s + '_gamma.png')
34 plt.show()
35 n = sum([i >= 25 for i in tmp])
36 print(1 + n / (sum([np.log(i / (25 - 0.5)) if i >= 25 else 0
   ↪ for i in tmp])))
37
38
39 def e2():
40     software = ['java', 'lucene']
41     for s in software:
42         generate_plot(s)
43         gamma(s)

```

3 Errors and attacks on the Internet

Figures 5 and 6 represent plots of the errors and attacks on the Internet and random graph, respectively.

As we can observe, random failures of nodes in the internet does not effect the connectivity of the network. However, when we disconnect certain nodes with the highest node degrees the internet quickly fails to deliver information across the network since this nodes are likely routers that connect large parts of the internet. On the other hand, random graph is robust to random failures and malicious attacks due to the fact that each node has a node degree similar

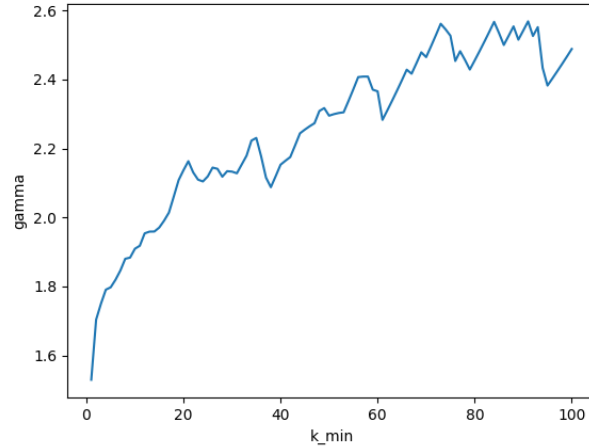


Figure 4: Lucene search engine.

to the average node degree.

The printout of the code used for this problem can be seen below.

```

1 def node_removal_list(n, p):
2     l = []
3     for i in n:
4         if np.random.choice(np.arange(0,2), p=[1-p, p]):
5             l.append(i)
6     return l
7
8
9 def e3():
10     fractions = [0, 0.1, 0.1, 0.1, 0.1, 0.1]
11     f = [0, .1, .2, .3, .4, .5]
12     pl = []
13     g = nx.read_adjlist('data/nec')
14     for p in fractions:
15         rem = node_removal_list(list(g.nodes), p)
16         g.remove_nodes_from(rem)
17         largest_cc = max(nx.connected_components(g), key=len)
18         pl.append(len(largest_cc) / g.number_of_nodes())
19     plt.plot(f, pl)
20     pl = []
21     g = nx.read_adjlist('data/nec')
22     for p in fractions:
23         rem = int(p*g.number_of_nodes())
24         l = list(g.degree())
25         l.sort(key=lambda x: x[1], reverse=True)

```

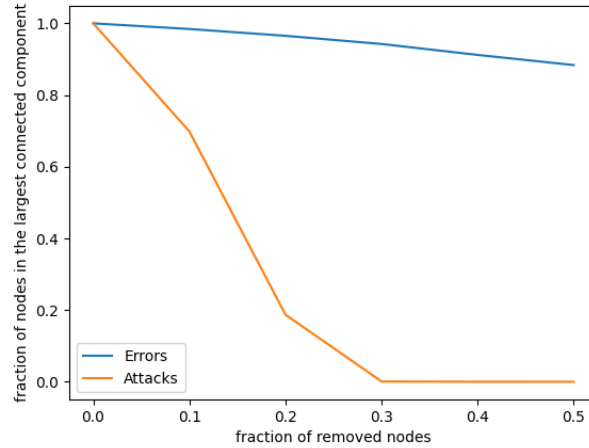


Figure 5: Internet.

```

26     tmp = [x for x, _ in l[:rem]]
27     g.remove_nodes_from(tmp)
28     largest_cc = max(nx.connected_components(g), key=len)
29     pl.append(len(largest_cc) / g.number_of_nodes())
30 plt.plot(f, pl)
31 plt.legend(['Errors', 'Attacks'], loc='lower left')
32 plt.xlabel('fraction of removed nodes')
33 plt.ylabel('fraction of nodes in the largest connected
    ↳ component')
34 plt.savefig('internet.png')
35 plt.show()
36
37 pl = []
38 g = nx.read_adjlist('data/nec')
39 n = g.number_of_nodes()
40 m = g.number_of_edges()
41 g = nx.gnm_random_graph(n, m)
42 for p in fractions:
43     rem = node_removal_list(list(g.nodes), p)
44     g.remove_nodes_from(rem)
45     largest_cc = max(nx.connected_components(g), key=len)
46     pl.append(len(largest_cc) / g.number_of_nodes())
47 plt.plot(f, pl)
48 pl = []
49 g = nx.gnm_random_graph(n, m)
50 for p in fractions:
51     rem = int(p * g.number_of_nodes())

```

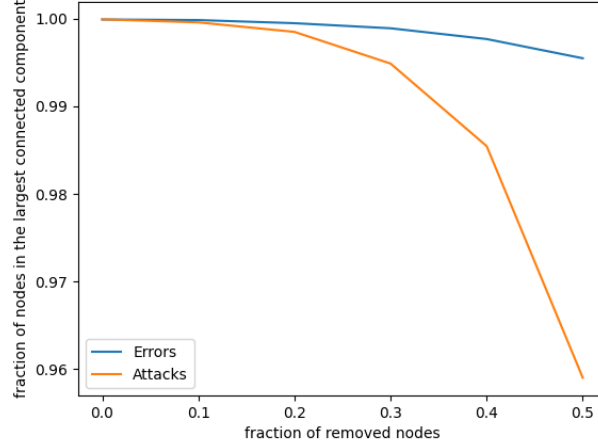


Figure 6: Random graph.

```

52     l = list(g.degree())
53     l.sort(key=lambda x: x[1], reverse=True)
54     tmp = [x for x, _ in l[:rem]]
55     if tmp:
56         g.remove_nodes_from(tmp)
57         largest_cc = max(nx.connected_components(g), key=len)
58         pl.append(len(largest_cc) / g.number_of_nodes())
59 plt.plot(f, pl)
60 plt.legend(['Errors', 'Attacks'], loc='lower left')
61 plt.xlabel('fraction of removed nodes')
62 plt.ylabel('fraction of nodes in the largest connected
    ↪ component')
63 plt.savefig('er.png')
64 plt.show()

```

4 HIV and network sampling

Firstly, scale-free networks have a power-law degree distribution. We can observe from the shape of the plots in the figure 7 that the original social network has a straight line shape which typical for power-law distribution and the sampled network has a bit curvy line thus it is not scale-free.

Secondly, small-world networks have high average clustering coefficient $\langle C \rangle$ and short average distance $\langle d \rangle \simeq \frac{\log n}{\log \langle k \rangle}$. Let us observe the average clustering coefficients and average distances of both network in the table 2.

As we can see, both network have a relatively short average distances. How-

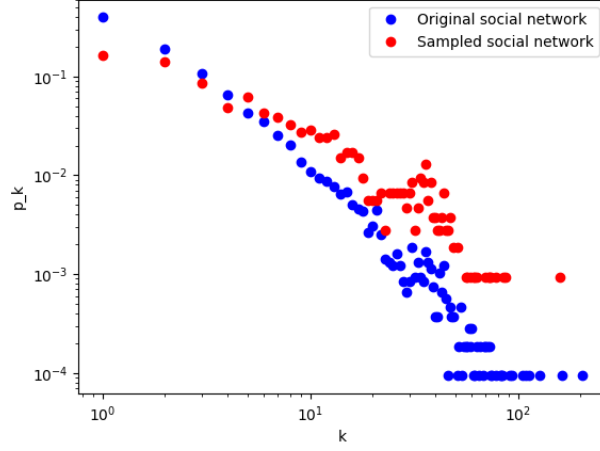


Figure 7: Node degree distributions of the original and sampled social network.

Network	$\langle C \rangle$	$\langle d \rangle$	$\frac{\log n}{\log \langle k \rangle}$
Original	0.27	7.49	6.12
Sampled	0.46	5.75	3.47

Table 2: Measures required for small-world network.

ever, the original network has a low clustering coefficient thus is not small-world network. On the other hand, the sampled network has high clustering therefore is a small-world network.

The printout of the code used for this problem can be seen below.

```

1 def distances(graph, i):
2     d = {i: np.infty for i in list(graph.nodes)}
3     q = deque()
4     q.append(i)
5     d[i] = 0
6     while q:
7         i = q.popleft()
8         for j in graph[i]:
9             if d[j] == np.infty:
10                 d[j] = d[i] + 1
11                 q.append(j)
12     return d
13
14
15 def average_distance(g):
16     d = list()

```



```

17     for i in list(g.nodes):
18         d.append(distances(g, i))
19     n = g.number_of_nodes()
20     return sum([v for tmp in d for v in tmp.values()]) / (n *
    ↪ (n-1))
21
22
23 def e4():
24     g = nx.read_adjlist('data/social')
25     rand_node =
    ↪ list(g.nodes)[np.random.randint(g.number_of_nodes())]
26
27     walk = {rand_node}
28     while True:
29         if len(walk) * 10 >= g.number_of_nodes():
30             break
31         adj = list(g.adj[rand_node])
32         rand_node = adj[np.random.randint(len(adj))]
33         walk.add(rand_node)
34     g_ind = g.subgraph(walk)
35     print('AIDS <C> = ' + str(nx.average_clustering(g_ind)))
36     print('AIDS <d> = ' + str(average_distance(g_ind)))
37     print('AIDS : ' + str(np.log(g_ind.number_of_nodes()) /
    ↪ np.log(sum([x for _, x in list(g_ind.degree())] /
    ↪ g_ind.number_of_nodes()))))
38     print('Social <C> = ' + str(nx.average_clustering(g)))
39     print('Social <d> = ' + str(average_distance(g)))
40     print('S : ' + str(np.log(g.number_of_nodes()) / np.log(sum([x
    ↪ for _, x in list(g.degree())] / g.number_of_nodes()))))
41
42     tmp = [x for _, x in g.degree()]
43     p_k = [tmp.count(i) / g.number_of_nodes() for i in
    ↪ range(max(tmp) + 1)]
44     plt.plot(p_k, 'bo')
45     tmp = [x for _, x in g_ind.degree()]
46     p_k = [tmp.count(i) / g_ind.number_of_nodes() for i in
    ↪ range(max(tmp) + 1)]
47     plt.plot(p_k, 'ro')
48     plt.yscale('log')
49     plt.xscale('log')
50     plt.legend(['Original social network', 'Sampled social
    ↪ network', 'Out-degree distribution'],
51               loc='upper right')
52     plt.xlabel('k')
53     plt.ylabel('p_k')
54     plt.savefig('aids.png')
55     plt.show()

```

5 Who to vaccinate?

Considering the friendship paradox which state that on average most people have fewer friend than their friends, the scheme where we select a number of individuals and than vaccinate a random acquaintance of theirs provide a better immunization.