

Homework #2

This homework is complete and will not be changed. The homework does not require a lot of writing, but may require a lot of thinking. It does not require a lot of processing power, but may require efficient programming. It accounts for 12.5% of the course grade. All questions and comments regarding the homework should be directed to [Piazza](#).

Submission details

This homework is due on **April 13th** at 2:00pm, while late days expire on **April 17th** at 1:00pm. The homework must be submitted as a hard-copy in the submission box in front of R 2.49 and also as an electronic version to [eUcilnica](#). It can be prepared in either English or Slovene and either written by hand or typed on a computer. The hard-copy should include (1) this cover sheet with filled out time of the submission and signed honor code, (2) short answers to the questions, which can also demand proofs, tables, plots, diagrams and other, and (3) a printout of all the code required to complete the exercises. The electronic submission should include only (1) answers to the questions in a single file and (2) all the code in a format of the specific programming language. Note that hard-copies will be graded, while electronic submissions will be used for plagiarism detection. The homework is considered submitted only when both versions have been submitted. Failing to include this honor code in the submission will result in **10% deduction**. Failing to submit all the developed code to [eUcilnica](#) will result in **50% deduction**.

Honor code

The students are strongly encouraged to discuss the homework with other classmates and form study groups. Yet, each student must then solve the homework by herself or himself without the help of others and should be able to redo the homework at a later time. In other words, the students are encouraged to collaborate, but should not copy from one another. Referring to any solutions obtained from classmates, course books, previous years, found online or other, is considered an honor code violation. Also, stating any part of the solutions in class or on [Piazza](#) is considered an honor code violation. Finally, failing to name the correct study group members, or filling out the wrong date or time of the submission, is also considered an honor code violation. Honor code violation will not be tolerated. Any student violating the honor code will be reported to **faculty disciplinary committee** and vice dean for education.

Name & SID: Jernej Vivod, 63160328

Study group: /

Date & time: 13.4.2020, 10:00

I acknowledge and accept the honor code.

Signature: 

Homework #2

Jernej Vivod
Introduction to Network Analysis

April 13, 2020

Question 1

Rank the importance of a dolphin named SN100 in the provided social network of bottlenose dolphins. Consider at least the suggested measures of node importance and state the rank of SN100 for each measure of node importance and briefly discuss the results.

We evaluated the estimated importance of dolphins using degree centrality, closeness centrality, betweenness centrality and PageRank score. The comparison of 10 highest-ranked dolphins and dolphin SN100 is visualized on figures 1 (degree centrality and closeness centrality) and 2 (betweenness centrality and PageRank score).

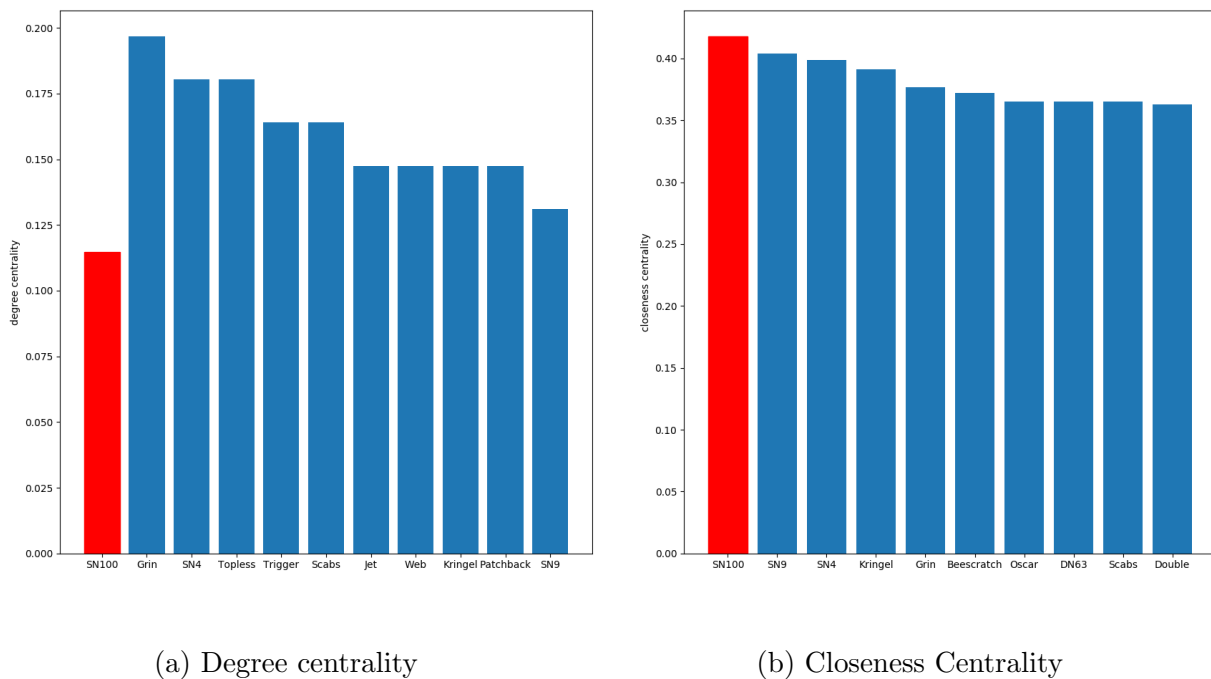


Figure 1: Degree centrality and closeness centrality of SN100 compared to other highest-ranked dolphins.

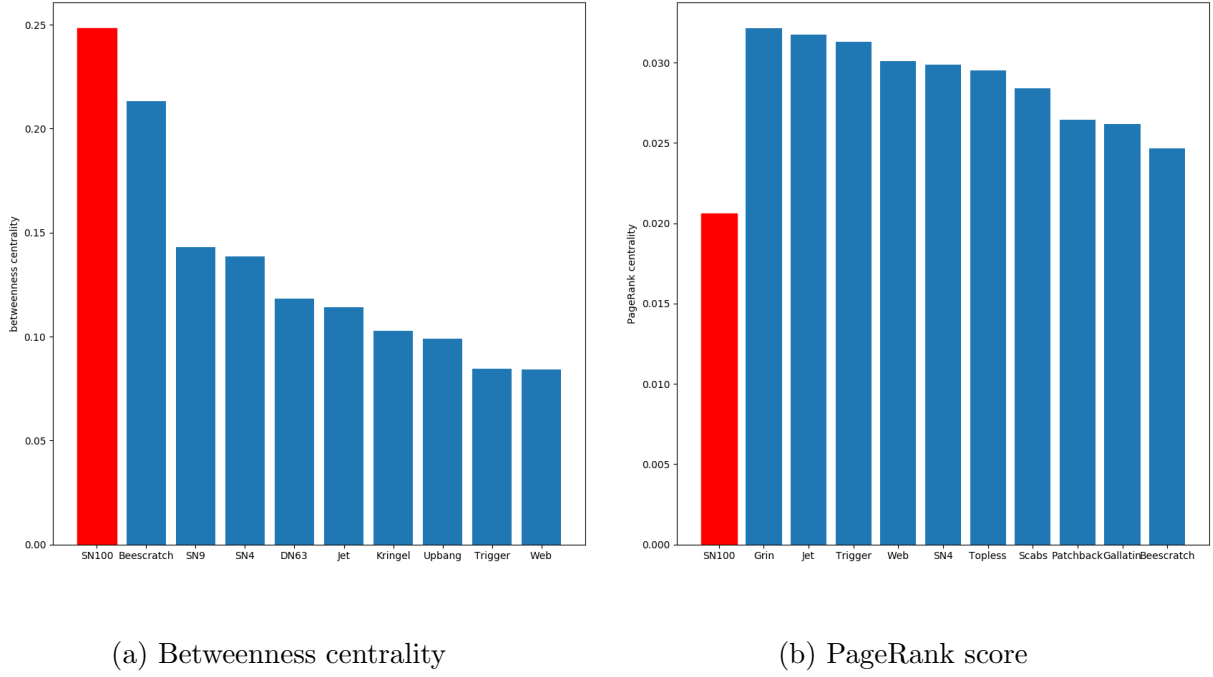


Figure 2: Betweenness centrality and PageRank centrality of SN100 compared to other highest-ranked dolphins.

	Rank	Percentile
Degree centrality	15	74.1935
Closeness centrality	1	100.0
Betweenness centrality	1	100.0
PageRank score	17	74.1935

Table 1

SN100 is ranked 15th for the degree centrality measure, 1st for the closeness centrality measure, 1st for the betweenness centrality measure and 17th for the PageRank measure.

Table 1 shows the rankings as well as the percentiles of the SN100's importance estimates. Thus the node representing the dolphin SN100 is estimated to be very close to other nodes and also to be part of a relatively large number of shortest paths between nodes.

The code used to provide the results is given below.

```

1 import networkx as nx
2 import parse_network
3
4
5 def node_importances(graph, measure):
6     """
7     Compute node importance scores according to specified measure.
8     Author: Jernej Vivod
9 
```

```

10     Args:
11         graph (obj): Networkx representation of the graph
12         measure (str): Argument specifying the node importance measure to use.
13     Returns:
14         (dict): Dictionary mapping node indices to their estimated importances.
15     """
16
17     # Check if specified measure valid.
18     if measure not in {'degree centrality', 'PageRank', 'betweenness', 'closeness'}:
19         raise(ValueError("the measure parameter can take the values of 'degree centrality',
20                             'PageRank', 'betweenness' or 'closeness'"))
21
22     # Compute node importances according to specified measure.
23     if measure == 'degree centrality':
24         return nx.degree_centrality(graph)
25     elif measure == 'PageRank':
26         return nx.pagerank(graph)
27     elif measure == 'betweenness':
28         return nx.betweenness_centrality(graph)
29     elif measure == 'closeness':
30         return nx.closeness_centrality(graph)
31
32 def node_rank(graph, idx_node, measure):
33     """
34     Compute node importance rank according to specified measure.
35     Author: Jernej Vivod
36
37     Args:
38         graph (obj): Networkx representation of the graph
39         idx_node (str): Index of node for which to compute the rank.
40         measure (str): Argument specifying the node importance measure to use.
41
42     Returns:
43         (int): Ranking of specified node according to specified importance measure.
44     """
45
46     # Check if specified measure valid.
47     if measure not in {'degree centrality', 'PageRank', 'betweenness', 'closeness'}:
48         raise(ValueError("the measure parameter can take the values of 'degree centrality',
49                             'PageRank', 'betweenness' or 'closeness'"))
50
51     # Compute rank of node according to specified measure.
52     if measure == 'degree centrality':
53         est_imp = nx.degree_centrality(graph)
54         return [key for key, val in sorted(est_imp.items(), key=lambda x: x[1], reverse=True)
55                 ].index(idx_node)
56     elif measure == 'PageRank':
57         est_imp = nx.pagerank(graph)
58         return [key for key, val in sorted(est_imp.items(), key=lambda x: x[1], reverse=True)
59                 ].index(idx_node)
60     elif measure == 'betweenness':
61         est_imp = nx.betweenness_centrality(graph)
62         return [key for key, val in sorted(est_imp.items(), key=lambda x: x[1], reverse=True)
63                 ].index(idx_node)
64     elif measure == 'closeness':
65         est_imp = nx.closeness_centrality(graph)
66         return [key for key, val in sorted(est_imp.items(), key=lambda x: x[1], reverse=True)
67                 ].index(idx_node)
68
69 def data_most_important(importance_dict, n_most_important, include_additional=None):
70     """
71     Compute data for making a bar plot of n nodes with highest estimated importance.
72     Author: Jernej Vivod
73
74     Args:
75         importance_dict (dict): Dictionary mapping node indices to their estimated

```

```

    importances
72     n_most_important (int): Number of most important nodes to keep
73     include_additional (str): Index of additional node from graph to include in plot
    data.
74     This node's data is appended to the front of the resulting lists.
75 Returns:
76     (tuple): tuple of lists of names of n most important nodes and importance scores of
    these nodes.
77 """
78
79 # Sort nodes by their evaluated importance.
80 sorted_nodes = [(k, v) for k, v in sorted(importance_dict.items(), key=lambda el: el[1],
    reverse=True)]
81 x = [nx.get_node_attributes(graph, 'name')[el[0]] for el in sorted_nodes[:n_most_
    important]]
82 y = [el[1] for el in sorted_nodes[:n_most_important]]
83
84 # If including additional specified node, add data for it if not yet present.
85 if include_additional and nx.get_node_attributes(graph, 'name')[include_additional] not
    in x:
86     x.insert(0, nx.get_node_attributes(graph, 'name')[include_additional])
87     y.insert(0, list(filter(lambda x: x[0] == include_additional, sorted_nodes))[0][1])
88
89 return x, y
90
91
92 if __name__ == '__main__':
93     import matplotlib.pyplot as plt
94     from scipy import stats
95
96     # Set name of dolphin of interest.
97     DOLPHIN_NAME = 'SN100'
98
99     # Parse the bottlenose dolphin network as well as associated data.
100    graph = parse_network.parse_network("../data/dolphins", create_using=nx.Graph)
101
102    # Get index of node corresponding to the dolphin of interest.
103    idx_dolphin = [el[0] for el in nx.get_node_attributes(graph, 'name').items() if el[1] ==
    'SN100'][0]
104
105    ### Bar charts of centralities ###
106    importances_degree centrality = node_importances(graph, 'degree centrality')
107    rank1 = node_rank(graph, idx_dolphin, 'degree centrality')
108    x_bar1, y_bar1 = data_most_important(importances_degree centrality, 10, include_
    additional=idx_dolphin)
109    fig1, ax1 = plt.subplots()
110    barlist1 = ax1.bar(x_bar1, y_bar1)
111    barlist1[0].set_color('r')
112    plt.ylabel("degree centrality")
113    print("degree centrality for dolphin {0}: {1:.4f}".format(DOLPHIN_NAME, y_bar1[0]))
114    print("degree centrality rank for dolphin {0}: {1}".format(DOLPHIN_NAME, rank1+1))
115    print("percentile of degree centrality for dolphin {0}: {1:.4f}".format(DOLPHIN_NAME,
    stats.percentileofscore(list(importances_degree centrality.values()), y_bar1[0])))
116
117    importances_pagerank = node_importances(graph, 'PageRank')
118    rank2 = node_rank(graph, idx_dolphin, 'PageRank')
119    x_bar2, y_bar2 = data_most_important(importances_pagerank, 10, include_additional=idx_
    dolphin)
120
121    fig2, ax2 = plt.subplots()
122    barlist2 = ax2.bar(x_bar2, y_bar2)
123    barlist2[0].set_color('r')
124    plt.ylabel("PageRank centrality")
125    print("PageRank centrality for dolphin {0}: {1:.4f}".format(DOLPHIN_NAME, y_bar2[0]))
126    print("PageRank centrality rank for dolphin {0}: {1}".format(DOLPHIN_NAME, rank2+1))
127    print("percentile of PageRank centrality for dolphin {0}: {1:.4f}".format(DOLPHIN_NAME,
    stats.percentileofscore(list(importances_pagerank.values()), y_bar2[0])))
128
129    importances_betweenness centrality = node_importances(graph, 'betweenness')
130

```

```

131 rank3 = node_rank(graph, idx_dolphin, 'betweenness')
132 x_bar3, y_bar3 = data_most_important(importances_betweenness centrality, 10, include_
    additional=idx_dolphin)
133 fig3, ax3 = plt.subplots()
134 barlist3 = ax3.bar(x_bar3, y_bar3)
135 barlist3[0].set_color('r')
136 plt.ylabel("betweenness centrality")
137 print("betweenness centrality for dolphin {0}: {1:.4f}".format(DOLPHIN_NAME, y_bar3[0]))
138 print("betweenness centrality rank for dolphin {0}: {1}".format(DOLPHIN_NAME, rank3+1))
139 print("percentile of betweenness centrality for dolphin {0}: {1:.4f}".format(DOLPHIN_
    NAME,
140     stats.percentileofscore(list(importances_betweenness centrality.values()), y_bar3
        [0])))
141
142 importances_closeness centrality = node_importances(graph, 'closeness')
143 rank4 = node_rank(graph, idx_dolphin, 'closeness')
144 x_bar4, y_bar4 = data_most_important(importances_closeness centrality, 10, include_
    additional=idx_dolphin)
145 fig4, ax4 = plt.subplots()
146 barlist4 = ax4.bar(x_bar4, y_bar4)
147 barlist4[0].set_color('r')
148 plt.ylabel("closeness centrality")
149 print("closeness centrality for dolphin {0}: {1:.4f}".format(DOLPHIN_NAME, y_bar4[0]))
150 print("closeness centrality rank for dolphin {0}: {1}".format(DOLPHIN_NAME, rank4+1))
151 print("percentile of closeness centrality for dolphin {0}: {1:.4f}".format(DOLPHIN_NAME,
152     stats.percentileofscore(list(importances_closeness centrality.values()), y_bar4[0]))
    )
153
154 plt.show()
155
156 #####

```

Question 2

Plot degree, in-degree and out-degree distributions on a doubly-logarithmic plot. Compare all three degree distributions p_k of each network and highlight the differences between them. Do the distributions appear to be scale-free?

Figure 3 shows the degree, in-degree and out-degree distributions for the *Java Namespace of Java Language* network plotted on a doubly-logarithmic plot. The dashed line represents the power-law model deduced by maximum-likelihood estimation. Figure 4 shows the same information for the *Lucene search engine library* network.

We can see that for both networks, only the in-degree distribution features the characteristic shape of a power-law distribution. Both the degree and out-degree distributions have a notable peak at a domain value some distance away from 0. This can be explained by observing the fact that very few classes implement such specific functionalities as to only be useful as a dependency for a very small set of other classes - hence the peak is located at considerable distance from 0. Similar reasoning can be applied to the degree distribution. However, most classes depend only on a small subset of other classes with a notable quantity of them depending on an unusually high number of classes. Many of such classes are classes that implement some main functionality of a library or program and therefore require many functionalities from more specialized classes. The deduced power-law exponents are stated by table 2. The code used to provide the results is given below.

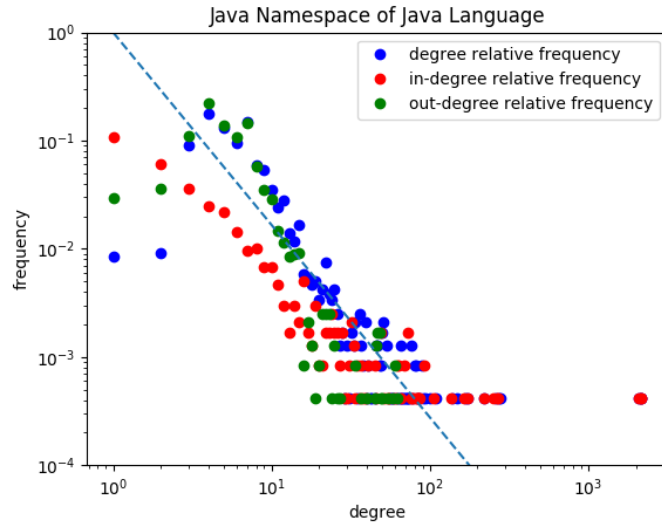


Figure 3: Degree, in-degree, out-degree and fitted power-law distribution model for the *Java Namespace of Java Language* network.

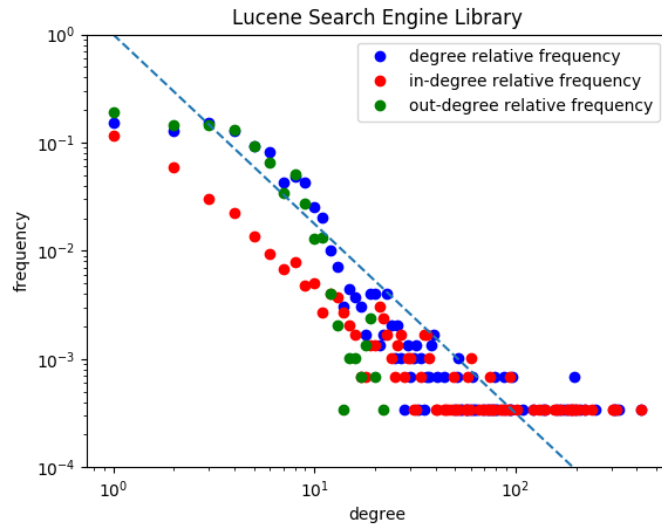


Figure 4: Degree, in-degree, out-degree and fitted power-law distribution model for the *Lucene search engine library* network.

	Power-law exponent
Java Namespace of Java Language	1.7789
Lucene Search Engine Library	1.7503

Table 2: Power-law exponents for the *Java Namespace of Java Language* and *Lucene Search engine* networks.

```
1 import networkx as nx
```

```

2 import math
3 import collections
4 import matplotlib.pyplot as plt
5
6
7 def plot_degree_distributions(graph):
8     """
9     Plot degree, in-degree and out-degree distribution on a doubly logarithmic plot.
10    Author: Jernej Vivod
11
12    Args:
13        graph (obj): Networkx graph representation.
14    """
15
16    # Compute relative degree, in-degree and out-degree frequencies.
17    degree_count = collections.Counter(dict(graph.degree()).values())
18    sum_deg_count = sum(degree_count.values())
19    degree_dist = {key: degree_count[key]/sum_deg_count for key in degree_count.keys()}
20    assert sum(degree_dist.values()) - 1.0 < 1.0e-4
21
22    in_degree_count = collections.Counter(dict(graph.in_degree()).values())
23    sum_in_deg_count = sum(in_degree_count.values())
24    in_degree_dist = {key: in_degree_count[key]/sum_in_deg_count for key in in_degree_count.
25                      keys()}
25    assert sum(in_degree_dist.values()) - 1.0 < 1.0e-4
26
27    out_degree_count = collections.Counter(dict(graph.out_degree()).values())
28    sum_out_deg_count = sum(out_degree_count.values())
29    out_degree_dist = {key: out_degree_count[key]/sum_out_deg_count for key in out_degree_
30                      count.keys()}
30    assert sum(out_degree_dist.values()) - 1.0 < 1.0e-4
31
32    # Plot relative degree frequencies on doubly-logarithmic plot.
33    fig, ax = plt.subplots()
34    ax.loglog(list(degree_dist.keys()), list(degree_dist.values()), 'bo', label="degree
35              relative frequency")
36    ax.loglog(list(in_degree_dist.keys()), list(in_degree_dist.values()), 'ro', label="in-
37              degree relative frequency")
38    ax.loglog(list(out_degree_dist.keys()), list(out_degree_dist.values()), 'go', label="out
39              -degree relative frequency")
40    ax.legend()
41
42    return fig, ax
43
44 def power_law_exponent(degrees, min_degree):
45     """
46     Evaluate power-law exponent gamma using maximum-likelihood estimate.
47     Author: Jernej Vivod
48
49     Args:
50         degrees (list): List of node degrees
51         min_degree (int): Degree cut-off
52
53     Returns:
54         (float): maximum-likelihood estimate of the power-law exponent.
55     """
56
57     # Count number of nodes falling below the threshold.
58     n = len(list(filter(lambda x: x >= min_degree, degrees)))
59
60     # Compute maximum-likelihood estimate for gamma.
61     return 1 + n*((sum([math.log(degree/(min_degree - 0.5)) for degree in degrees if degree
62                       >= min_degree]))*(-1.0))
63
64 if __name__ == '__main__':

```



```

64     ### Parse graphs ###
65     PATH1 = '../data/java'
66     PATH2 = '../data/lucene'
67     graph_java = nx.read_edgelist(PATH1, create_using=nx.DiGraph)
68     graph_lucene = nx.read_edgelist(PATH2, create_using=nx.DiGraph)
69
70     # Set plot titles and axis labels.
71     title1 = "Java Namespace of Java Language"
72     title2 = "Lucene Search Engine Library"
73     xlabel = "degree"
74     ylabel = "frequency"
75
76     # Plot exponent estimate or not.
77     PLOT_EXPONENT_ESTIMATE = True
78
79     # Plot degree, in-degree and out-degree distributions.
80     fig1, ax1 = plot_degree_distributions(graph_java)
81     fig2, ax2 = plot_degree_distributions(graph_lucene)
82
83     # Set titles and axis labels.
84     ax1.set_title(title1)
85     ax2.set_title(title2)
86     ax1.set_xlabel(xlabel)
87     ax1.set_ylabel(ylabel)
88     ax2.set_xlabel(xlabel)
89     ax2.set_ylabel(ylabel)
90
91     # Compute power-law exponent using maximum-likelihood estimation.
92     gamma1 = power_law_exponent(dict(graph_java.in_degree()).values(), min_degree=3)
93     gamma2 = power_law_exponent(dict(graph_lucene.in_degree()).values(), min_degree=3)
94
95     # If plotting exponent estimate on plot.
96     if PLOT_EXPONENT_ESTIMATE:
97         dom1 = range(1, max(dict(graph_java.in_degree()).values())+1)
98         dom2 = range(1, max(dict(graph_lucene.in_degree()).values())+1)
99         ax1.loglog(dom1, [el**(-gamma1) for el in dom1], '--')
100        ax2.loglog(dom2, [el**(-gamma2) for el in dom2], '--')
101        ylim = (1.0e-4, 1)
102        ax1.set_ylim(ylim[0], ylim[1])
103        ax2.set_ylim(ylim[0], ylim[1])
104
105    # Show plot.
106    plt.show()

```

Question 3

Consider the provided Internet overlay map represented by an undirected graph. Study how removing a fraction of random nodes affects the fraction of nodes in the largest connected component. Study how removing a fraction of nodes with the highest degrees affects the fraction of nodes in the largest connected component. Repeat the same computations on an Erdős–Rényi random graph model with parameters n and m taken from the Internet overlay map network. Briefly discuss the results.

Figure 5 shows the fraction of nodes remaining in the largest connected component after a fraction of either randomly selected nodes or nodes with highest degrees (hubs) were removed. We can see that removing nodes of very high degree quickly deteriorates the connectedness of the *Internet overlay map* network. This is due to these nodes acting as hubs

and forming vital bridges within the network. oppositely, removing fractions of randomly chosen nodes in the *Internet overlay map* network does not have a large detrimental effect on the connectedness of the network as the probability of choosing a vital hub that is important in ensuring the network stays connected is relatively low.

On the other hand, the Erdős–Rényi model does not feature prominent hubs or nodes of high bridgeness and the fraction of nodes in the largest connected component is therefore very robust to removal of random as well as high-degree nodes.

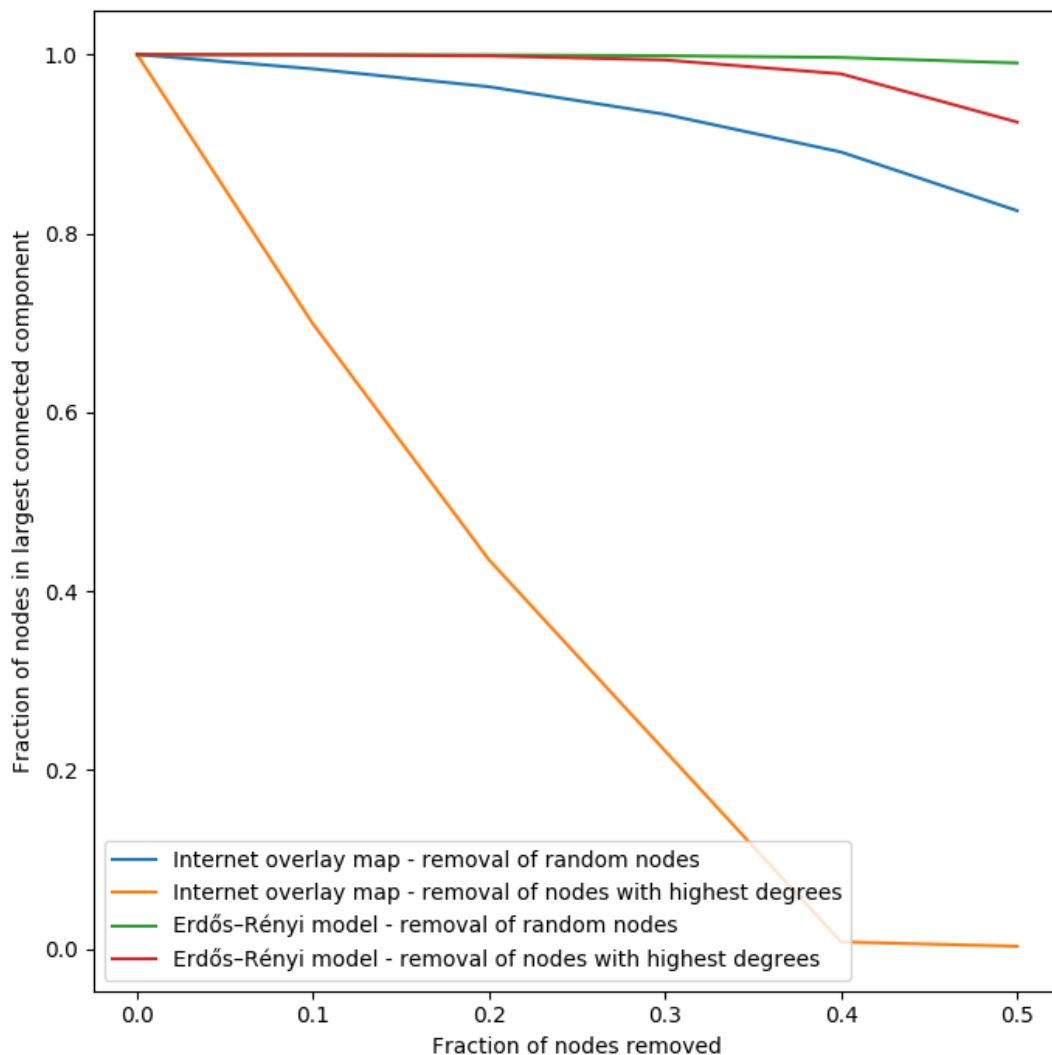


Figure 5: Fraction of nodes in the largest connected component after a fraction of either randomly selected nodes or nodes with highest degrees (hubs) were removed.

The code used to provide the results is given below.

```
1 import networkx as nx
2 import parse_network
3 import random
4
5 def remove_frac_nodes(graph, frac, remove_hubs):
6     """
7     Remove a fraction of nodes in specified graph.
8     Author: Jernej Vivod
9
10    Args:
11        graph (obj): Networkx representation of a graph.
12        frac (float): The fraction of nodes to remove from the graph.
13        remove_hubs (bool): If set to true, remove specified fraction of
14        nodes with highest degree. Else select nodes to be removed randomly.
15
16    Returns:
17        (obj): Networkx representation of a graph with specified
18        fraction of nodes removed.
19    """
20
21    if frac < 0.0 or frac > 1.0:
22        raise ValueError("Fraction must be between 0.0 and 1.0")
23    else:
24        if remove_hubs:
25            # Remove fraction of nodes with highest degree.
26            to_remove = [key for (key, val) in sorted(dict(graph.degree()).items(), key=
27                        lambda x: x[1], reverse=True))[:round(frac*graph.number_of_nodes())]
28            graph.remove_nodes_from(to_remove)
29            return graph
30        else:
31            # Remove fraction of randomly selected nodes from graph.
32            graph.remove_nodes_from(random.sample(graph.nodes(), round(frac*graph.number_of_
33                        nodes()))))
34            return graph
35
36 def components(graph):
37     """
38     Find connected component in undirected graph.
39     Author: Jernej Vivod
40
41    Args:
42        graph (obj): Networkx representation of a network
43
44    Returns:
45        (list): List of lists containing node IDs representing connected components.
46    """
47
48    # Empty list for storing the connected components
49    connected_components = []
50
51    # Perform DFS to find connected components.
52    while graph.number_of_nodes() > 0:
53        connected_components.append(component(graph))
54
55    # Return connected components.
56    return connected_components
57
58 def component(graph):
59     """
60     Find next connected component in graph.
61     Author: Jernej Vivod
62
63    Args:
```

```

64         graph (obj): Networkx representation of a network
65
66     Returns:
67         (list): the list containing the node IDs constituting the connected component.
68     """
69
70     # Empty list for storing the next connected component
71     component = []
72
73     # Stack for implementing DFS
74     stack = []
75
76     # Add starting node to stack.
77     root_node = list(graph.nodes())[0]
78     stack.append(root_node)
79
80     # Perform DFS.
81     while len(stack) > 0:
82
83         # Pop next node from stack.
84         node_nxt = stack.pop()
85
86         # If node in graph.
87         if node_nxt in graph:
88
89             # Append index of current node to component list.
90             component.append(node_nxt)
91
92             # Go over neighbors of current node.
93             for el in graph.neighbors(node_nxt):
94
95                 # If neighbor not yet visited, add to stack.
96                 if el in graph:
97                     stack.append(el)
98
99             # Remove current node from graph.
100             graph.remove_node(node_nxt)
101
102     # Return found connected component.
103     return component
104
105
106 def frac_in_lcc(graph):
107     """
108     Compute fraction of nodes in largest connected component.
109     Author: Jernej Vivod
110
111     Args:
112         graph (obj): Networkx representation of a graph.
113
114     Returns:
115         (float): Fraction of nodes in largest connected component.
116     """
117
118     cc = components(graph.copy())
119     return max(map(lambda x: len(x), cc))/graph.number_of_nodes()
120
121
122 if __name__ == '__main__':
123     import matplotlib.pyplot as plt
124
125     # Parse network.
126     PATH = "../data/nec"
127     graph = parse_network.parse_network(PATH, create_using=nx.Graph)
128
129     # Construct Erdos-Renyi model with same number of nodes and edges.
130     graph_er_model = nx.gnm_random_graph(graph.number_of_nodes(), graph.number_of_edges())
131

```

```

132 # Initialize list of fractions of nodes to remove.
133 fracs = [0, 0.1, 0.2, 0.3, 0.4, 0.5]
134
135 # Initialize lists for storing results.
136 frac_lcc_graph_rm_rand = []
137 frac_lcc_graph_rm_hubs = []
138 frac_lcc_graph_er_rm_rand = []
139 frac_lcc_graph_er_rm_hubs = []
140
141 # Go over list of fractions.
142 for frac in fracs:
143     print(frac)
144
145     # Remove fraction of nodes from graph and Erdos-Renyi model.
146     graph_rm_rand = remove_frac_nodes(graph.copy(), frac, False)
147     graph_rm_hubs = remove_frac_nodes(graph.copy(), frac, True)
148     graph_er_rm_rand = remove_frac_nodes(graph_er_model.copy(), frac, False)
149     graph_er_rm_hubs = remove_frac_nodes(graph_er_model.copy(), frac, True)
150
151     # Add computed fractions of nodes in largest connected component to results lists.
152     frac_lcc_graph_rm_rand.append(frac_in_lcc(graph_rm_rand))
153     frac_lcc_graph_rm_hubs.append(frac_in_lcc(graph_rm_hubs))
154     frac_lcc_graph_er_rm_rand.append(frac_in_lcc(graph_er_rm_rand))
155     frac_lcc_graph_er_rm_hubs.append(frac_in_lcc(graph_er_rm_hubs))
156
157 # Plot fractions of nodes in largest connected component with respect to fraction of
158 # nodes removed.
159 fig, ax = plt.subplots()
160 ax.plot(frac, frac_lcc_graph_rm_rand, label="Internet overlay map - removal of random
161 nodes")
162 ax.plot(frac, frac_lcc_graph_rm_hubs, label="Internet overlay map - removal of nodes
163 with highest degrees")
164 ax.plot(frac, frac_lcc_graph_er_rm_rand, label="Erdos-Renyi model - removal of random
165 nodes")
166 ax.plot(frac, frac_lcc_graph_er_rm_hubs, label="Erdos-Renyi model - removal of nodes
167 with highest degrees")
168 ax.legend()
169 plt.xlabel("Fraction of nodes removed")
170 plt.ylabel("Fraction of nodes in largest connected component")
171 plt.show()

```

Question 4

Perform a random walk on the provided social network until you sample 10% of the nodes and take the induced graph constructed from such sampled nodes. Is the original social network small-world and/or seemingly scale-free? Is the sampled network small-world and/or seemingly scale-free? Reason why and support your answers with appropriate computations.

Figure 6 shows the node degree distribution for original graph and sampled graph. We can see that the distribution in the original plot follows a power-law distribution and the sampled graph also follows an approximate power-law distribution. When the fraction of nodes in the sample decreases, the frequency of low degree nodes decreases somewhat as there is a higher chance of the random walk sampling nodes with a higher degree and we notice a peak in the distribution somewhat distant from 0.

Table 3 shows the average shortest distance and average clustering of the original and sampled graphs. We can notice that the shortest paths are short especially when compared

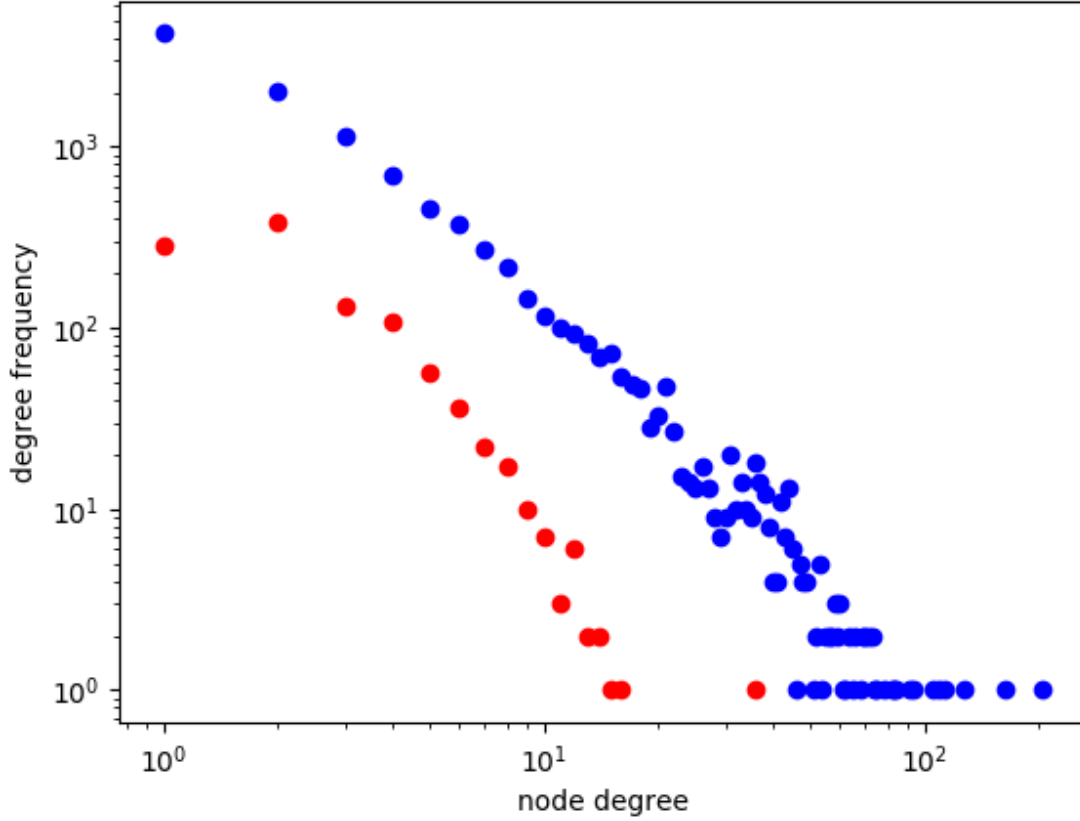


Figure 6: Degree distribution in original graph and in graph induced by random walk.

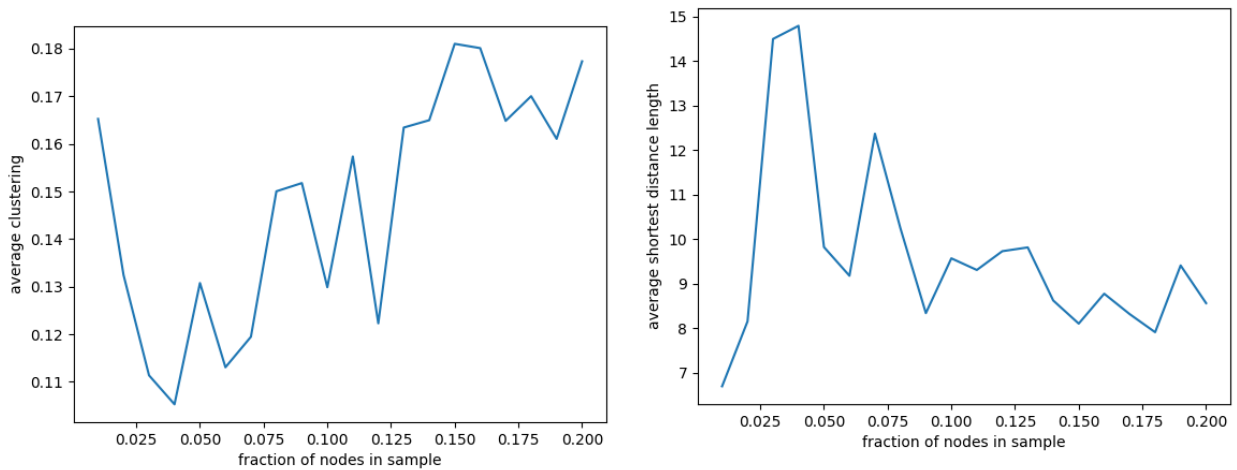
to the large number of nodes constituting the network. The clustering coefficient is also much higher than expected by chance. Both networks can therefore also be classified as small-world.

	Average shortest distance	Average clustering
Original graph	7.4855	0.2659
Sampled graph (10%)	8.5413	0.1353

Table 3: Average shortest path distance and average clustering for original and sampled graph.

It is also interesting to plot the average clustering and average shortest path lengths for different sample sizes as shown on figure 7. We can see that the average shortest path length soon settles in close vicinity of the average shortest path length in the original graph. The average clustering is somewhat lower but still much higher than would be expected by chance.

The code used to provide the results is given below.



(a) Average clustering for varying sample sizes. (b) Average shortest distance for varying sample sizes.

Figure 7: Average clustering and average shortest path length for varying sample sizes.

```

1 import networkx as nx
2 import parse_network
3 import random
4 import collections
5
6 def random_walk(graph, frac_sample=0.1):
7     """
8     Perform random walk on specified graph until specified fraction
9     of nodes have been covered. Return graph induced by such a random
10    walk.
11
12    Args:
13        (obj): Networkx graph representation.
14        frac_sample: Fraction of nodes in graph to cover.
15
16    Returns:
17        (obj): Networkx graph representation of the graph induced
18        by the random walk.
19    """
20
21    # Initialize set of visited nodes.
22    visited = set()
23
24    # Initialize set of traversed edges.
25    walk_edges = set()
26
27    # Randomly choose starting node.
28    start_node = random.choice(list(graph.nodes()))
29    node_current = start_node
30
31    # While specified fraction of network not covered, perform random walk.
32    while len(visited) < frac_sample*graph.number_of_nodes():
33
34        # Visit randomly chosen neighbor of current node and add to results sets.
35        node_nxt = random.choice(list(graph.neighbors(node_current)))
36        visited.add(node_nxt)
37        walk_edges.add((node_current, node_nxt))
38        node_current = node_nxt

```

```

39
40     # Construct graph from traversed edges and return it.
41     ind_graph = nx.Graph()
42     ind_graph.add_edges_from(walk_edges)
43     return ind_graph
44
45
46
47
48 if __name__ == '__main__':
49     import matplotlib.pyplot as plt
50
51     # Parse network.
52     PATH = "../data/social"
53     graph = parse_network.parse_network(PATH, create_using=nx.Graph)
54
55     # Get graph induced by random walk that covers 10% of the nodes.
56     ind_graph = random_walk(graph, 0.06)
57
58     # Print average distance and clustering in original and sampled graph.
59     # print("Average distance - original graph: {0}".format(nx.average_shortest_path_length(
60     graph)))
61     print("Average distance - sampled graph: {0}".format(nx.average_shortest_path_length(ind
62     _graph)))
63
64     # print("Average clustering - original graph: {0}".format(nx.average_clustering(graph)))
65     print("Average clustering - sampled graph: {0}".format(nx.average_clustering(ind_graph))
66     )
67
68     # Plot degree distributions for original and sampled graphs.
69     degree_freq_original = collections.Counter(dict(graph.degree()).values())
70     x_original, y_original = list(degree_freq_original.keys()), list(degree_freq_original.
71     values())
72
73     degree_freq_induced = collections.Counter(dict(ind_graph.degree()).values())
74     x_ind, y_ind = list(degree_freq_induced.keys()), list(degree_freq_induced.values())
75
76     fig, ax = plt.subplots()
77     ax.loglog(x_original, y_original, 'bo')
78     ax.loglog(x_ind, y_ind, 'ro')
79     plt.xlabel("node degree")
80     plt.ylabel("degree frequency")
81     plt.show()

```

Question 5

Probability that an individual would spread the disease through their social network is proportional to k^2 , where k is the degree of the corresponding node. Consider two immunization schemes for preventing the spread of diseases. In the first scheme, you randomly select some number of individuals and vaccinate them. In the second scheme, you randomly select the same number of individuals, but then rather vaccinate a random acquaintance of theirs. Which of the two schemes do you expect to provide better immunization? Why?

According to the so-called friendship paradox, which can be explained as a form of sampling bias, the friends or acquaintances of a person have, on average, a greater number of friends or acquaintances than that particular person. By this fact, the second scheme should be more successful as there is a higher chance of selecting individuals with a large social

circle and therefore considerably higher probability of spreading the disease through their social network.

We can simulate this empirically by writing a short script that simulates the procedure some number of times and reports the normalized sum of squared degrees of nodes representing unvaccinated people as a measure of the effectiveness of the vaccination program.

Table 4 shows the results for 30 repetitions of both vaccination schemes where 10% of subjects are randomly selected.

	Measure
First scheme	77.3959
Second scheme	53.2773

Table 4: Sum of squared degrees of nodes representing unvaccinated persons for both vaccination schemes.

We can see a notable difference in the normalized sum of squared nodes representing unvaccinated persons that is stable when repetitions and fractions of selected sample are altered. As predicted by computations supporting the so-called friendship paradox, the second vaccination scheme is notably better in targeting individuals with a high probability of spreading the disease through their social circle.

The code used to simulate the vaccination schemes is given below.

```

1 import networkx as nx
2 import parse_network
3 import random
4
5
6 def mark_random_nodes(graph, num_to_mark):
7     """
8     Mark random specified number of nodes in specified graph by returning
9     them in a list of node IDs.
10
11     Args:
12         graph (obj): Networkx representation of a graph.
13         num_to_mark (int): number of nodes to mark (return in a list).
14
15     Returns:
16         (list): List of marked nodes.
17     """
18
19     # Randomly sample specified number of nodes.
20     return random.sample(list(graph.nodes()), num_to_mark)
21
22
23 def mark_random_nodes_neighbors(graph, num_to_mark):
24     """
25     Select specified number of random nodes and then randomly mark/select
26     a neighbor of each one and return a list of such marked nodes.
27
28     Args:
29         graph (obj): Networkx representation of a graph.
30         num_to_mark (int): number of nodes to mark (return in a list).
31
32     Returns:
33         (list): List of marked nodes.
34     """
35

```

```

36     # Initialize list for storing marked nodes.
37     marked = []
38
39     # Randomly sample specified number of nodes.
40     sample = random.sample(list(graph.nodes()), num_to_mark)
41
42     # Go over sampled nodes and mark randomly chosen neighbors.
43     for node in sample:
44         marked.append(random.choice(list(graph.neighbors(node))))
45
46     # Return list of marked nodes.
47     return marked
48
49
50 if __name__ == '__main__':
51
52     # Parse network.
53     PATH = '../data/social'
54     graph = parse_network.parse_network(PATH, create_using=nx.Graph)
55
56     # Fraction of nodes to mark and number of runs to perform.
57     FRAC_TO_MARK = 0.1
58     NUM_RUNS = 30
59
60     # Initialize results aggregate.
61     aggr1 = 0
62     aggr2 = 0
63
64     for idx in range(NUM_RUNS):
65
66         # Randomly mark specified fraction of nodes.
67         marked_scheme1 = mark_random_nodes(graph, int(round(graph.number_of_nodes()*FRAC_TO_MARK)))
68
69         # Randomly mark random neighbors of specified fraction of nodes.
70         marked_scheme2 = mark_random_nodes_neighbors(graph, int(round(graph.number_of_nodes()*FRAC_TO_MARK)))
71
72         # Sum squared degrees of unmarked nodes normalized by number of nodes in graph and add to aggregate.
73         sum_norm_squared1 = sum([graph.degree()[e1]**2/graph.number_of_nodes() for e1 in graph.nodes() if e1 not in marked_scheme1])
74         sum_norm_squared2 = sum([graph.degree()[e1]**2/graph.number_of_nodes() for e1 in graph.nodes() if e1 not in marked_scheme2])
75         aggr1 += sum_norm_squared1
76         aggr2 += sum_norm_squared2
77
78     # Compute average sum of normalized squared node degrees for the two marking schemes.
79     avg_sum_norm_squared1 = aggr1/NUM_RUNS
80     avg_sum_norm_squared2 = aggr2/NUM_RUNS
81
82     # Print results.
83     print("Average sum of normalized squared node degrees for first marking scheme ({0} marked, {1} runs): {2:.4f}".format(NUM_RUNS, FRAC_TO_MARK, avg_sum_norm_squared1))
84     print("Average sum of normalized squared node degrees for second marking scheme ({0} marked, {1} runs): {2:.4f}".format(NUM_RUNS, FRAC_TO_MARK, avg_sum_norm_squared2))

```