

Homework #1

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 March 23rd at 2:00pm, while late days expire on March 27th 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

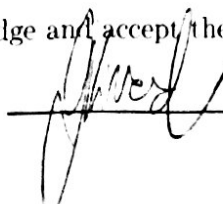
The students are strongly encouraged to discuss the homework with other classmates and form study groups. Yet, each student must then solve the homework by herself or himself without the help of others and should be able to redo the homework at a later time. In other words, the students are encouraged 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: VERNEJ VIVOD

Study group: 1

Date & time: 22.3.2020, 13:50

I acknowledge and accept the honor code.

Signature: 

Homework #1

Jernej Vivod
Introduction to Network Analysis

March 15, 2020

Question 1.1

What can you say about the degree sequences $\{k\}$ and degree distributions p_k of the shown graphs?

The second graph follows a power law distribution as there exist some nodes with an unusual number of connections (hubs). The first network has nodes that are generally well-connected. There is an absence of hubs. Plotting the degree distribution would produce a characteristic bell-shaped curve (binomial distribution). Unlike the second network, this network contains some nodes with degree 1.

Question 1.2

Assuming a simple undirected network with n nodes, m links and c connected components. Show that the given two inequalities hold. Using these inequalities give a criterion for m that ensures a connected network. Is the criterion practically useful? Why?

Let us start the proof by proving that if $n - c \leq m$ holds for a graph with n nodes, then it also holds for a graph with $n + 1$ nodes. Note that the upper bound for the number of links in a fully connected network with n nodes is $m_{max} = \frac{n(n-1)}{2}$. Using this, we can show:

$$\begin{aligned}(n+1) - c &\leq \frac{(n+1)n}{2} \\ 2n+2-2c &\leq n^2+n \\ n+2-2c &\leq n^2 \\ -2c+2 &\leq n^2-n\end{aligned}$$

Since the right side of the inequality is strictly positive for $n > 0$ and the left side of the inequality is non-positive for $c > 0$, the equality holds for all graphs with more than one node and more than one connected component.

We now only need to show, that the equality holds for a base case. Assume a graph with a single node (and one connected component). We can see that:

$$\begin{aligned}
n - c &\leq m \\
1 - 1 &\leq 0 \\
0 &\leq 0
\end{aligned}$$

By the induction step, we now know that the inequality holds for any graph with $n > 0$ nodes.

The reasoning for the second inequality can be given as follows. Suppose we have a network with n nodes and c connected components. The upper bound for the number of links in such a network is achieved by having the network split into one giant fully connected component and $c - 1$ other "trivial" connected components consisting of single nodes with no links. We know that the number of links in a fully connected network is given by the number of unique pairs of nodes in the network $\binom{n}{2}$. By having the aforementioned network configuration, we "lose" $c - 1$ nodes to the trivial connected components. By this reasoning, $m \leq \binom{n-c+1}{2}$ must hold.

To ensure a connected network, the inequality $m \geq n - 1$ must hold. To link n nodes into a connected network, we need to place exactly one link between a node with no links and another node that is part of a connected component. To create such a tree graph, we need to place $n - 1$ links.

This criterion is useful as it gives the minimum number of connections needed to form a path from each node to every other node which is useful in many real-world applications (certain infrastructure projects, etc.).

Question 1.3

What would the mentioned algorithm for finding connected components of undirected networks find in a directed network if one could follow the links in any direction? What would it find in a directed network if one could follow the links only in the proper direction? What would it find in a directed network if one could follow the links only in the opposite direction? Design an efficient algorithm for finding strongly connected components in directed networks.

If one could follow the links in any direction, the algorithm would find the weakly connected components of the directed graph. If one could follow the links only in the proper direction, the algorithm would find all the nodes reachable from the starting node by following the directed paths in their proper directions as well as the nodes reachable from the nodes not in the component of the initial node but without the nodes in the previously found components. Similarly, the algorithm would find the same types of components in the transpose graph (with the edge directions reversed). In the original graph, that would be equivalent to finding the nodes from which the to starting node can be reached by following the directed edges.

The algorithm for finding the strongly connected components can be formulated by observing the mentioned properties as well as the behaviour of depth-first search. We can

design an algorithm that will find all nodes reachable from root node as well as to find the subset of the found nodes that can reach the root node by using depth-first search on the transpose graph.

Running the algorithm we see that the Enron e-mail communication network has 78058 strongly connected components. The size of the largest one is 9164. There are relatively many (especially single node) strongly connected components in the network (32.15% the number of nodes) This is due to a large number of people that only received e-mail communication but did not reply. The people who did reply overwhelmingly became part of the largest strongly connected component. The results are surprising in terms of the large number of single-node strongly connected components (one-way communication).

The pseudocode for the algorithm used to find the strongly connected components is given below:

```

Data: Directed graph
Result: List of strongly connected components
Initialize stack S ;
Put root node on stack S ;
while S not empty do
    Pop node from stack S. ;
    if popped node not yet visited then
        if node has unvisited neighbors then
            push popped node back on stack S. ;
            Add unvisited neighbors to stack S. ;
        else
            Add node on DFS finish stack. ;
            Pop node from stack. ;
        end
    else
        Add node on DFS finish stack. ;
        Pop node from stack. ;
    end
end
Get transpose graph G' ;
Initialize empty list SCC ;
while DFS finish stack not empty do
    Pop node from DFS finish stack ;
    Perform DFS starting from popped node to get next SCC and add it to SCC ;
end

```

Algorithm 1: Algorithm used to find the strongly connected components.

The implementation of the algorithm used to find the strongly connected is given below:

```

1 import networkx as nx
2
3
4 def strongly_connected_components(graph):
5     """
6     Find strongly connected components in directed graph.
7     Author: Jernej Vivod (vivod.jernej@gmail.com)
8
9     Args:
10         graph (networkx.classes.digraph.DiGraph): directed graph.
11
12     Returns:
13         (list): List of lists containing nodes in connected components.
14     """
15
16 def get_finish_stack(graph):
17     """
18     Get stack of nodes based on their DFS finish times.
19
20     Author:
21         Jernej Vivod (vivod.jernej@gmail.com)
22
23     Args:
24         graph (networkx.classes.digraph.DiGraph): directed graph.
25
26     Returns:
27         finish_stack (list): stack of nodes ordered by their DFS finish times.
28     """
29
30     # Initialize starting finish time enumeration value.
31     step_nxt = 1
32
33     # Initialize dictionary for storing start and finish times.
34     finish_dict = {el:[0, 0] for el in graph.nodes()}
35
36     # Initialize stack for storing the nodes in finish time order.
37     finish_stack = []
38
39     # Enumerate while there are nodes in the graph.
40     while graph.number_of_nodes() > 0:
41         step_nxt = dfs_enumerate(graph, finish_dict, finish_stack, step_nxt)
42
43     # Return stack of nodes based on their finish time.
44     return finish_stack
45
46
47 def dfs_enumerate(graph, finish_dict, finish_stack, step):
48     """
49     Enumerate component based on DFS finish times (auxiliary function)
50     Author: Jernej Vivod (vivod.jernej@gmail.com)
51
52     Args:
53         graph (networkx.classes.digraph.DiGraph): directed graph.
54         finish_dict (dict): dictionary mapping nodes to their DFS start and finish times
55
56         finish_stack (list): stack for keeping nodes based on their DFS finish times.
57         step (int): DFS enumeration start value.
58
59     Returns:
60         (int): next value to use in enumeration of nodes based on DFS finish times.
61     """
62
63     # Initialize stack for performing DFS.
64     stack = []
65
66     # Add starting node to stack.
67     node_start = list(graph.nodes())[0]

```

```

67     stack.append(node_start)
68
69     # While stack not empty, perform DFS and enumerate nodes
70     # based on finish time.
71     while len(stack) > 0:
72
73         # Get node on top of stack.
74         node_current = stack[-1]
75
76         # If node in graph ...
77         if graph.has_node(node_current):
78
79             # Set start time and increment enumeratio value.
80             finish_dict[node_current][0] = step
81             step += 1
82
83             # Get neighbors and remove node from graph.
84             neighbors = graph.neighbors(node_current)
85             graph.remove_node(node_current)
86             has_unvisited_neighbors = False
87
88             # Go over neighbors and add unvisited to stack.
89             for neighbor in neighbors:
90                 if graph.has_node(neighbor):
91                     has_unvisited_neighbors = True
92                     stack.append(neighbor)
93
94             # If node has no unvisited neighbors ...
95             if not has_unvisited_neighbors:
96
97                 # If not yet finished, set finish time and pop from stack.
98                 if finish_dict[node_current][1] == 0:
99                     finish_dict[node_current][1] = step
100                     finish_stack.append(node_current)
101                     step += 1
102                     stack.pop()
103             else:
104                 # If not yet finished, set finish time and pop from stack.
105                 if finish_dict[node_current][1] == 0:
106                     finish_dict[node_current][1] = step
107                     finish_stack.append(node_current)
108                     step += 1
109                     stack.pop()
110
111     # Return last enumeration value.
112     return step
113
114
115 def get_strongly_connected_components(graph_trans, finish_stack):
116     """
117     Get strongly connected components by performing DFS on transposed graph
118     and utilizing the computed stack of nodes ordered by their DFS finish
119     times on the original graph.
120     Author: Jernej Vivod (vivod.jernej@gmail.com)
121
122     Args:
123         graph_trans (networkx.classes.digraph.DiGraph): The transpose of the original
124         graph.
125         finish_stack (list): stack of nodes ordered by their DFS finish times.
126
127     Returns:
128         (list): list of lists of nodes constituting the strongly connected components.
129     """
130
131     # Initialize list for sotring connected components.
132     components = []
133
134     # Initialize stack for performing DFS.

```

```

134     stack = []
135
136     # While stack of nodes ordered by finish time not empty ...
137     while len(finish_stack) > 0:
138
139         # Add starting node for DFS to stack.
140         node_start = finish_stack.pop()
141         stack.append(node_start)
142
143         # Initialize list for storing the next
144         # strongly connected component.
145         component = []
146
147         # While stack not empty, perform DFS.
148         while len(stack) > 0:
149             current_node = stack.pop()
150
151             # If node in graph ...
152             if graph_trans.has_node(current_node):
153
154                 # Add node to component.
155                 component.append(current_node)
156
157                 # Get neighbors and remove node from graph.
158                 neighbors = graph_trans.neighbors(current_node)
159                 graph_trans.remove_node(current_node)
160
161                 # Go over neighbors and add unvisited to stack.
162                 for neighbor in neighbors:
163                     if graph_trans.has_node(neighbor):
164                         stack.append(neighbor)
165
166             # If found non-empty component, add to list of
167             # strongly connected components.
168             if len(component) > 0:
169                 components.append(component)
170
171         # Return list of strongly connected components.
172         return components
173
174
175     # Get stack of nodes ordered by their DFS finish times.
176     finish_stack = get_finish_stack(graph.copy())
177
178     # Get transpose of graph.
179     graph_trans = graph.reverse()
180
181     # Compute strongly connected components.
182     return get_strongly_connected_components(graph_trans, finish_stack)

```

Question 1.4

Think of another example of a network for which $\langle C \rangle \rightarrow \text{const}$ and $C \rightarrow 0$ when $n \rightarrow \infty$.

Another example of a network with such properties would be the "petal" network shown on figure [1](#). To prove that the network has such properties, let us denote with k the number of the petals in the network. In a network with 2 petals, there are 4 connected triples that extend between the petals as well as 3 connected triples on each on the petals themselves. Now consider a network with k petals. To find the number of connected triplets that extend

between pairs of petals, we need to count such triplets in $\frac{n(n-1)}{2}$ unique pairs of petals. We also need to count the $3 \cdot k$ triplets found in each individual petal. As the number of petals increases, the growth of the term capturing the increase in the number of connected triplets is clearly quadratic as more and more petals are added, whereas the number of triangles in the graph increases by 1 for each added petal (linear growth). The clustering coefficient of each node in the petals is equal to 1. As more and more petals are added, the average clustering converges to the constant value of 1.

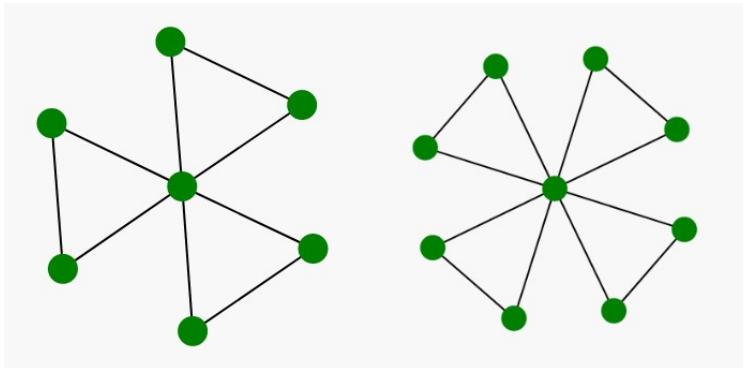


Figure 1: Examples of a petal network with 3 and 4 petals.

Question 1.5

Design an efficient algorithm for computing E_{90} of a connected undirected network. Implement the algorithm and compute E_{90} for citation networks of physics papers. Are the results surprising? Why? Compute also the number of nodes and the average degree $\langle k \rangle$ of all three networks and discuss the results.

We can formulate the algorithm by using two similar but distinct definitions of the effective diameter. We can either compute the lengths of paths between every unique pair of connected vertices or take the average minimal distance required to reach 90% of nodes from each node. Both methods were implemented for this task. Table 1 summarizes the results.

	all-pairs E_{90}	unique-pairs E_{90}	number of nodes	average degree $\langle k \rangle$
Citations for 2010-2011	14.3792	15.0	18985	4.3694
Citations for 2010-2012	11.9346	13.0	38356	5.7302
Citations for 2010-2013	10.4929	12.0	56473	7.0955

Table 1: 90-percentile effective diameter, number of nodes and average degree for citations networks of physics papers for increasing time intervals.

We can see that analyzing citations over a larger interval produces a lower 90-percentile effective diameter and a higher average node degree. The decrease in the 90-percentile effective diameter can seem surprising at first but can be interpreted by observing that taking a larger time interval increases the chances of early papers being cited by papers

towards the end of the time interval which in turn increases the average node degree. These citations by later papers often then form links that reduce distances between earlier papers as well as later ones and thus decrease the 90-percentile effective diameter of the network.

The pseudocode below presents the algorithm used to solve this task using distances between unique pairs of nodes. To compute the distances using all pairs of nodes, a slightly modified algorithm which computes the 90th percentile for distances between sampled node and all other nodes in the graph was used. The latter algorithm obtains the results by averaging the 90th percentile distances from each node to every other node in the graph.

Data: Graph G

Result: 90-percentile effective diameter of graph G

idx = 1 ;

for *each node in sorted nodes (by ID) in G* **do**

 dists := get next distances to sorted G. ;

 /* select relevant distances (not yet computed for that pair of nodes). */

 dists := dists[idx:end] ;

 idx := idx + 1 ;

end

return value at 90th percentile in dists. ;

Algorithm 2: Algorithm used to find the 90-percentile effective diameter.

The code written to find the solutions to this question is given below:

```

1 import networkx as nx
2 import numpy as np
3
4 def effective_diameter(graph, mode, percentile):
5     """
6     Compute nth-percentile effective diameter of graph.
7     Author: Jernej Vivod (vivod.jernej@gmail.com)
8
9     Args:
10         graph (networkx.classes.digraph.DiGraph): Graph for which to compute the nth-
11             percentile effective diameter.
12         mode (str): Method of computing the results. If equal to 'unique_pairs', compute
13             result as nth-percentile of distances between unique node pairs.
14         percentile (int): Percentile used in the computations.
15
16     Returns:
17         (float): The nth-percentile effective diameter of specified graph.
18     """
19
20     def pairwise_distances(graph):
21         """
22         Compute vector of distances between unique pairs of nodes in graph.
23
24         Author:
25             Jernej Vivod (vivod.jernej@gmail.com)
26
27         Args:
28             graph (networkx.classes.digraph.DiGraph): Graph for which to compute the nth-
29                 percentile effective diameter.
30
31         Returns:

```

```

30         (numpy.ndarray): Vector of distances between unique pairs of nodes in specified
31         graph.
32
33     """
34     # Allocate vector for storing unique pairwise distances.
35     num_nodes = graph.number_of_nodes()
36     dists = np.empty(int((num_nodes*(num_nodes-1))/2), dtype=int)
37
38     # Set start position for selecting relevant pairwise distances.
39     # Set index for relevant pairwise distances vector.
40     start_pos = 1
41     dists_idx = 0
42
43     # Go over nodes and compute pairwise distances.
44     for node in sorted(map(int, graph.nodes())):
45
46         # Get distances from next node to all other nodes.
47         dists_nxt = get_distances(graph, str(node))
48
49         # Get relevant pairwise distances
50         dists[dists_idx:dists_idx+len(dists_nxt[start_pos:])] = dists_nxt[start_pos:]
51         start_pos += 1
52         dists_idx += len(dists_nxt[start_pos:]) + 1
53     return dists
54
55
56 def distances_percentile(graph, percentile):
57     """
58     Compute vector of nth-percentile distances to every node from each node.
59
60     Author:
61         Jernej Vivod (vivod.jernej@gmail.com)
62
63     Args:
64         graph (networkx.classes.digraph.DiGraph): Graph for which to compute the nth-
65             percentile effective diameter.
66         percentile (int): Percentile used in the computations.
67
68     Returns:
69         (numpy.ndarray): Vector of n-th percentile distances to every node from each node
70         .
71     """
72
73     # Allocate vector for storing 90th percentile distances.
74     num_nodes = graph.number_of_nodes()
75     dists_perc = np.empty(num_nodes, dtype=float)
76
77     # Go over nodes and compute distances at percentiles.
78     for (idx, node) in enumerate(sorted(map(int, graph.nodes()))):
79
80         # Get distances from next node to all other nodes.
81         dists_nxt = get_distances(graph, str(node))
82
83         # Compute distance representing the percentile.
84         dists_perc[idx] = np.percentile(dists_nxt, percentile)
85
86     # Return vector of distances at percentiles for each node.
87     return dists_perc
88
89 def get_distances(graph, node):
90     """
91     Compute distances from node to every other node in graph.
92
93     Args:
94         graph (networkx.classes.digraph.DiGraph): Graph for which to compute the nth-

```

```

95         percentile effective diameter.
96         node (str): Node for which to compute distances to every other node.
97
98     Returns:
99         (numpy.ndarray): array of distances from specified node to every other node in
100         the graph.
101
102     Author:
103         Jernej Vivod (vivod.jernej@gmail.com)
104
105     """
106
107     # Initialize array for storing distances.
108     dists = np.full(graph.number_of_nodes(), -1, dtype=int)
109
110     # Set distance of current node to itself to zero.
111     dists[int(node)-1] = 0
112
113     # Initialize queue and add starting node.
114     queue = []
115     queue.append(node)
116
117     # While queue not empty, perform BFS.
118     while len(queue) > 0:
119         node_current = queue.pop(0)
120         for neighbor in graph.neighbors(node_current):
121
122             # Compute distances to neighbors.
123             if dists[int(neighbor)-1] == -1:
124                 dists[int(neighbor)-1] = dists[int(node_current)-1] + 1
125                 queue.append(neighbor)
126
127     # Return array of distances of node to all the other nodes.
128     return dists
129
130 # Compute nth-percentile effective diameter.
131 if mode == 'unique_pairs':
132     dists_vec = pairwise_distances(graph)
133     return np.percentile(dists_vec, percentile)
134 elif mode == 'all_pairs':
135     dists_perc = distances_percentile(graph, percentile)
136     return np.mean(dists_perc)

```

Question 2.1

Design an algorithm that does not select nodes uniformly at random, but proportional to their degrees. Thus, node i is selected with probability $\frac{k_i}{2m}$, where k_i is its degree and m is the number of links. The algorithm should run in constant time $O(n)$, whereas you can assume any standard network representation.

Since the probabilities are expressed as fractions, we do not need to resort to algorithms such as the alias method for sampling from weighted distributions. Instead, we can represent the network as a list of nodes (their IDs) where each node is repeated k_i times. The probability of a node being selected by random uniform sampling from such a list will be equal to the desired probability. The sampling can then simply be done arbitrary many times by generating a random integer in the interval $[1, \text{list.length}]$ which can be done in constant time. The pseudocode for the algorithm as well as the code used for the implementation are given below.

Data: Graph G

Result: Node selected with probability proportional to its degree.

$L :=$ list of node ID's where the i -th node ID is repeated k_i times. ;

/* Selection can now be done in $O(1)$ (constant) time. */

for *number of nodes to select* **do**

 idx := random integer from $[0, \text{length}(L)]$;
 yield $L[\text{idx}]$;

end

Algorithm 3: Algorithm used to select N nodes with probability proportional to their degree in constant time.

```
1 import random
2
3 def select_preferential(node_to_degree, n):
4     """
5     Select n nodes in specified graph with probability proportional
6     to their degrees (no replacement).
7     Author: Jernej Vivod
8
9     Args:
10        node_to_degree (dict): dictionary mapping node IDs to their degrees.
11        n (int): number of nodes to select.
12
13     Returns:
14        (list): IDs of selected nodes.
15     """
16
17     # Create list of node IDs where number of occurrences is equal to the node's degree.
18     sel_list = [key for key in node_to_degree.keys() for _ in range(node_to_degree[key])]
19
20     # Allocate list for storing results.
21     res = ['']*n
22
23     # Sample n nodes without replacement.
24     for idx_sel in range(n):
25         sel = sel_list[random.randint(0, len(sel_list)-1)]
26         res[idx_sel] = sel
27         sel_list = list(filter(lambda x: x != sel, sel_list))
28
29     # Return sample as list of IDs.
30     return res
```

Question 2.2

Consider a random graph model in which links are placed independently between each pair of nodes i and j with probability p_{ij} proportional to $v_i v_j$, where v_i is some non-negative number associated with node i . First show that the expected node degree k_i is proportional to v_i . Next, derive an exact expression for p_{ij} in terms of the degree sequence $\{k_i\}$ and discuss the result.

The expected node degree $E(k_i)$ can be written as:

$$\begin{aligned} E(k_i) &= \sum_{j \neq i} p_{ij} \\ E(k_i) &= \alpha \sum_{j \neq i} v_i v_j \\ E(k_i) &= \alpha v_i \sum_{j \neq i} v_j \end{aligned}$$

The last line clearly shows that the expected degree is proportional to the value of v_i .

The probability of a link between node i and node j can be expressed by dividing the number of ways to choose one half-edge from node i and one half-edge from node j by the number of ways to connect a half-edge to another half-edge. We can write this as:

$$p_{ij} = \frac{k_i k_j}{2m - 1}$$

We know that p_{ij} is proportional to $v_i v_j$. We can therefore write:

$$\begin{aligned} p_{ij} &= \frac{k_i k_j}{2m - 1} \\ v_i v_j &\propto \frac{k_i k_j}{2m - 1} \\ v_i v_j &\propto \frac{k_i k_j}{2m - 1} \\ v_i v_j &\propto k_i k_j \\ v_i v_j &= \alpha k_i k_j \end{aligned}$$

This shows that the product of the values $v_i v_j$ is directly proportional to the product of the degrees of nodes i and j .

Question 2.3

Represent a small part of the Facebook social network as an undirected graph and compute its degree distribution p_k . Construct an Erdős–Rényi random graph with parameters n and m and again compute p_k . Compute also the theoretical degree distribution of the Erdős–Rényi random graph (Poisson distribution). Finally construct a random graph according to the preferential attachment model using the specified procedure. Again compute p_k . Plot the results on a doubly logarithmic (log-log) plot. Compare all four degree distributions p_k and highlight the differences among them.

The results obtained by randomly sampling 30000 nodes from the Facebook social network graph are shown on figure [2](#).

We can see that both the Facebook social network network as well as the Barabási–Albert preferential attachment model follow a power-law node degree distribution, whereas the

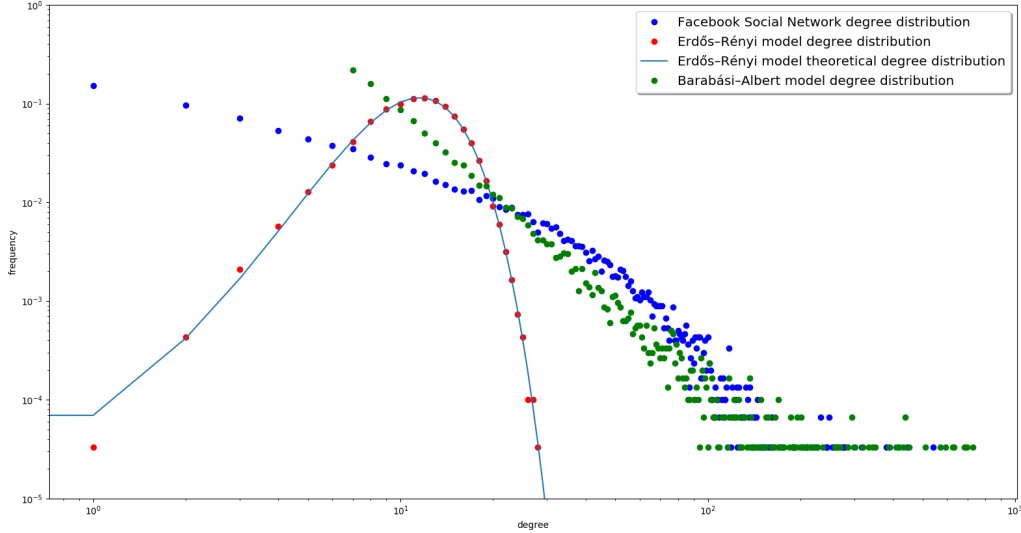


Figure 2: Degree distributions of a sample of the Facebook social network graph and its different models.

Erdős–Rényi model follows the binomial degree distribution which sharply contrasts with the more natural Barabási–Albert model that features hubs as well as a multitude of nodes with relatively low degrees. We can also note that the actual, empirically compute degree distribution of the Erdős–Rényi model closely follows the theoretically predicted distribution.

The code written to find the solutions to this question is given below:

```

1 import math
2 import random
3 import networkx as nx
4 from collections import Counter
5 import matplotlib.pyplot as plt
6
7 from select_preferential import select_preferential
8
9 # Parse Facebook Social Network graph from file.
10 GRAPH_PATH = '../data/facebook'
11 graph_fb_full = nx.read_edgelist(GRAPH_PATH, create_using=nx.Graph)
12
13 # Create a subgraph of the Facebook Social Network graph by randomly sampling nodes.
14 SAMPLE_SIZE_FB = 30000
15 sampled_nodes = random.sample(graph_fb_full.nodes, SAMPLE_SIZE_FB)
16 graph_fb = graph_fb_full.subgraph(sampled_nodes)
17
18 # Compute histogram of node degrees.
19 hist_fb = Counter([degree for n, degree in graph_fb.degree()])
20
21 # Compute mean node degree.
22 mean_k = sum([degree*hist_fb[degree]/graph_fb.number_of_nodes() for degree in hist_fb.keys()
23              ])
24
25 # Plot degree distribution for Facebook Social Network graph.
26 degrees = hist_fb.keys()
27 plt.loglog(list(degrees), [hist_fb[degree]/graph_fb.number_of_nodes() for degree in degrees

```

```

    ],
    'bo', label="Facebook Social Network degree distribution")
27
28
29 # Create Erdos-Renyi model with same number of nodes and links.
30 graph_model_er = nx.gnm_random_graph(graph_fb.number_of_nodes(), graph_fb.number_of_edges())
31
32 # Compute histogram of node degrees.
33 hist_model_er = Counter([degree for n, degree in graph_model_er.degree()])
34
35 # Plot degree distribution for model
36 degrees_model_er = hist_model_er.keys()
37 plt.loglog(list(degrees_model_er), [hist_model_er[degree]/graph_model_er.number_of_nodes()
    for degree in degrees_model_er],
38            'ro', label="Erdos-Renyi model degree distribution")
39
40
41 # Plot theoretical distribution of the Erdos-Renyi model degree distribution.
42 degrees_model_er_sorted = list(range(71))
43 hist_model_er_theoretical = [(mean_k**k)*math.exp(-mean_k)/math.factorial(k) for k in
    degrees_model_er_sorted]
44 plt.loglog(degrees_model_er_sorted, hist_model_er_theoretical, '-', label="Erdos-Renyi model
    theoretical degree distribution")
45
46
47 # Construct preferential attachment model (Barabasi-Albert model).
48
49 # Start with fully connected graph.
50 graph_model_ba = nx.complete_graph(math.ceil(mean_k)+1, create_using=nx.Graph)
51 node_nxt = max(list(graph_model_ba.nodes())) + 1
52
53 for idx in range(SAMPLE_SIZE_FB - math.ceil(mean_k) - 1):
54
55     graph_model_ba.add_node(node_nxt)
56
57     # Select ceil(mean_k/2) existing nodes with probability proportional to their degrees
    and link to them.
58     for node in select_preferential(dict(graph_model_ba.degree()), n=math.ceil(mean_k/2)):
59         graph_model_ba.add_edge(node_nxt, node)
60
61     node_nxt += 1
62     print("done {0}/{1}".format(idx, SAMPLE_SIZE_FB-math.ceil(mean_k)-2))
63
64 # Compute histogram of node degrees.
65 hist_model_ba = Counter([degree for n, degree in graph_model_ba.degree()])
66
67 # Plot degree distribution for model.
68 degrees_model_ba = hist_model_ba.keys()
69 plt.loglog(list(degrees_model_ba), [hist_model_ba[degree]/graph_model_ba.number_of_nodes()
    for degree in degrees_model_ba],
70            'go', label="Barabasi-Albert model degree distribution")
71
72
73 # Finish plotting.
74 plt.legend(loc='upper right', shadow=True, fontsize='x-large')
75 plt.xlabel("degree")
76 plt.ylabel("frequency")
77 plt.ylim(1e-5, 1)
78 plt.show()

```

Question 3

You are given the Slovenian highway network from 2010 with traffic loads at each location. Find out which measure of node position could be utilized to

best predict the traffic loads. You should consider at least node degree k_i , node clustering coefficient $c_i = \frac{2t_i}{k_i(k_i-1)}$ and the node harmonic mean distance $\ell_i^{-1} = \frac{1}{n-1} \sum_j \frac{1}{d_{ij}}$. Compute the Pearson or Spearman correlation coefficient between the values returned by some node measure and the actual traffic loads. Compute the correlation coefficient for each of the three measures. Are the results expected? Why? List also the top ten locations according to the best node measure along with the computed values and the actual traffic loads.

Table 2 shows the values of the Pearson correlation coefficient for each of the three measures. Table 3 lists the top ten locations according to the best node measure along with

Node measure	Pearson correlation coefficient value
node degree measure	0.2757
node clustering measure	undefined (0 clustering for all nodes)
harmonic mean distance measure	0.6178

Table 2: Values of the Pearson correlation coefficient for used node measures.

the computed values and the actual traffic loads.

Node measure	Pearson correlation coefficient value	Traffic load
Kozarje	0.1230	35759.1561
Koseze	0.1222	50232.0277
Zadobrova	0.1200	33779.1600
Malence	0.1195	55297.8647
Brdo	0.1183	35759.1561
Slivnica	0.1155	20031.1456
Vič	0.1137	20686.5315
Celovška	0.1133	28813.9702
Zaloška	0.1132	20686.5315

Table 3: Top ten locations according to the best node measure (harmonic mean distance) and the actual traffic load.

The predictions made by using the harmonic mean distance measure (closeness centrality of a node measure by the reciprocal of farness) are relatively good. We can see that the top rated nodes are located close to important junctions in the ring around Ljubljana which itself serves as a sort of a junction between the extending branches. The results are therefore expected.

The code written to find the solutions to this question is given below:

```

1 import networkx as nx
2 import scipy.stats as sps
3 import re
4
5 def load_with_attributes(path):
6     """
7     Load graph from specified file and add listed node attributes to graph.

```



```

8     Author: Jernej Vivod (vivod.jernej@gmail.com)
9
10    Args:
11        path (str): Path to file containing graph data.
12
13    Returns:
14
15
16
17    """
18
19    # Parse graph.
20    graph = nx.read_edgelist(GRAPH_PATH, create_using=nx.Graph)
21
22    # Initialize dictionaries for parsing data.
23    load = dict()
24    names = dict()
25
26    # Flag indicating whether next line contains data to be parsed.
27    parse_data = False
28    with open(path, 'r') as f:
29        for line in f:
30            if line[0] == "#":
31                if not parse_data:
32                    if len(line.split(" ")) == 1:
33                        parse_data = True
34                    elif parse_data and len(line.split(" ")) > 1:
35                        data_raw_nxt = line.split(" ")
36                        names[data_raw_nxt[1]] = re.findall(r'"([^\"]*)"', line)[0]
37                        load[data_raw_nxt[1]] = float(data_raw_nxt[-1].strip())
38            else:
39                nx.set_node_attributes(graph, load, 'load')
40                nx.set_node_attributes(graph, names, 'name')
41        return graph
42
43
44    # Parse Slovenian highways network dataset.
45    GRAPH_PATH = '../data/highways'
46    graph = load_with_attributes(GRAPH_PATH)
47
48    # Compute node degree, node clustering coefficient and the node harmonic mean distance.
49    degrees = graph.degree()
50    clustering_coefficients = nx.clustering(graph)
51    harmonic_mean_distances = {node : sum(map(lambda x: 1.0/x, filter(lambda x: x != 0,
52        nx.single_source_shortest_path_length(graph, node).values())))/(graph.number_of_nodes()
53        -1) for node in graph.nodes()}
54
55    # Compute correlation of measure values with actual loads.
56    nodes = graph.nodes()
57    nodes_to_data = graph.nodes(data=True)
58    loads_data = [nodes_to_data[node]['load'] for node in nodes]
59
60    # Compute pearson correlation coefficient for the node degree measure.
61    degree_measure_data = [degrees[node] for node in nodes]
62    degree_measure_pearson = sps.pearsonr(degree_measure_data, loads_data)
63    print("Pearson correlation coefficient for the node degree measure: {0}, two-sided p-value:
64        {1}".format(*degree_measure_pearson))
65
66    # Compute pearson correlation coefficient for the node clustering measure.
67    clustering_measure_data = [clustering_coefficients[node] for node in nodes]
68    clustering_measure_pearson = sps.pearsonr(clustering_measure_data, loads_data)
69    print("Pearson correlation coefficient for the node clustering measure: {0}, two-sided p-
70        value: {1}".format(*clustering_measure_pearson))
71
72    # Compute pearson correlation coefficient for the harmonic mean distance measure.
73    harmonic_mean_distances_measure_data = [harmonic_mean_distances[node] for node in nodes]
74    harmonic_mean_distances_measure_pearson = sps.pearsonr(harmonic_mean_distances_measure_data,

```

```

loads_data)
73 print("Pearson correlation coefficient for the harmonic mean distance measure: {0}, two-
    sided p-value: {1}".format(*harmonic_mean_distances_measure_pearson))
74
75
76 # List top 10 nodes according to the best measure. List the value of the computed measure
    and the actual load.
77
78 # Select best measure and get top rated nodes.
79 best_measure = max(((degree_measure_pearson[0], degrees, 'node degree'), (clustering_measure
    _pearson[0], clustering_coefficients, 'node clustering'),
80     (harmonic_mean_distances_measure_pearson[0], harmonic_mean_distances, 'node harmonic
        mean distance'))))
81 top_nodes = list(map(lambda x: x[0], sorted(best_measure[1].items(), key=lambda x: x[1],
    reverse=True)))[:10]
82
83 # Print results.
84 print("Best measure: '{0}'".format(best_measure[2]))
85 print("Top-rated nodes:")
86 for idx, node in enumerate(top_nodes):
87     print("{0} {1}: measure={2:.4f}, load={3:.4f}".format(idx+1, nodes_to_data[node]['name
        '], best_measure[1][node], nodes_to_data[node]['load']]))

```