

Homework #3

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 **May 4th** at 2:00pm, while late days expire on **May 8th** 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: 5.5.2020, 11:00

I acknowledge and accept the honor code.

Signature: 

Homework #3

Jernej Vivod
Introduction to Network Analysis

May 5, 2020

Question 1

Show that $L = B^T B$ (L is the graph Laplacian matrix and B is the oriented link incidence matrix). Using this equality further show that all eigenvalues of L are non-negative and that the vector of all ones is an eigenvector of L .

Consider an incidence matrix of a graph where the first node is either the first end-point or the second end point of all links going to other nodes connected with the first node. To recover the node degree from an oriented incidence matrix, we need to sum the absolute values of each column. This is equivalent to taking the dot product of each column with itself. Following the definition of matrix multiplication we can easily see that the elements on the diagonal of the product correspond to the degrees of the nodes.

$$(B^T B)_{i,i} = \sum_j B_{ij}^T B_{ji}$$

On the other hand, taking a dot product of a column with a column corresponding to a different node simply counts the links between the two nodes. Since the first end-point is positive and the second negative, the result of the dot product is equal to the negative sum of the number of links between the two nodes. In the case of a simple graph with no multiple edges, the result will simply be equal to -1 if the corresponding nodes are connected.

Combining these two observations, we can see, that the matrix product $B^T B$ results in a matrix with node degrees on the diagonal and negative values (-1 in the case of a simple graph) on in the i, j -th position, if the j -th and i -th nodes are connected. This can be written as

$$B^T B = D - A,$$

where D is the matrix with node degrees on the diagonal and A is the adjacency matrix. This is equal to the definition of the graph Laplacian matrix L .

A real $n \times n$ symmetric matrix whose eigenvalues are all non-negative is called a positive semi-definite matrix. We know that a matrix $M \in \mathbb{R}^{n \times n}$ is positive semi-definite if $\forall x \in \mathbb{R}^n, x^T M x \geq 0$.

We can easily show that the Laplacian matrix is positive semi-definite by using its incidence matrix decomposition we proved earlier.

$$x^T L x = x^T B^T B x = (Bx)^T Bx = \|Bx\|_2^2 \geq 0$$

Because each row of the graph Laplacian matrix contains the degree of its corresponding node on the diagonal position and negative values on positions corresponding to nodes to which the node corresponding to this node is linked, the values sum up to 0. Multiplying the graph Laplacian matrix with a vector of all ones corresponds to performing a row-wise sum. Therefore, the result will be a vector of all zeros. Thus, a vector of all ones is an eigenvector of the graph Laplacian matrix with an eigenvalue 0.

Question 2

Compute modularity Q of described partition and express it in terms of n_c and n . Find the size of clusters n_c that optimizes modularity Q and express it in terms of n .

We can easily write the terms in the formula for computing the modularity as functions of n and n_c . We know that the number of links in a community C in such a graph is $n_c - 1$ where n_c is the number of nodes in the community. We also know that the number of nodes is equal to the number of links. and that the total degree is simply two times the number of nodes in each partition.

We can rewrite the equation for computing the modularity as:

$$\begin{aligned} Q &= \sum_C \frac{m_c}{m} - \left(\frac{k_c}{2m}\right)^2 = \frac{n}{n_c} \left(\frac{n_c - 1}{n} - \left(\frac{n_c - 1}{2n}\right)^2\right) \\ &= \left(\frac{n_c - 1}{n} - \frac{4n_c^2}{4n}\right) \frac{n}{n_c} = \frac{nn_c - n - nn_c^2}{nn_c} = 1 - \frac{1}{n_c} - \frac{n_c}{n} \end{aligned}$$

We can optimize the modularity with respect to n_c by the standard procedure of differentiating the expression with respect to n_c , setting the derivative equal to 0 and solving for n_c .

$$\begin{aligned} \frac{\partial Q}{\partial n_c} &= 0 - \frac{1}{n_c^2} - \frac{1}{n} = 0 \\ \frac{1}{n_c^2} &= -\frac{1}{n} \\ n_c &= \sqrt{n} \end{aligned}$$

We can see that the modularity is optimized when we set the cluster size n_c equal to the square root of the number of nodes in the graph.

Question 3

Compare community detection methods using specified approach.

Figure 1 shows the visualization of a Girvan-Newman graph with three communities containing 24 nodes each, with expected degree equal to 20 and with the μ parameter equal to 0.1. The code describing the implementation is given at the end of this section.

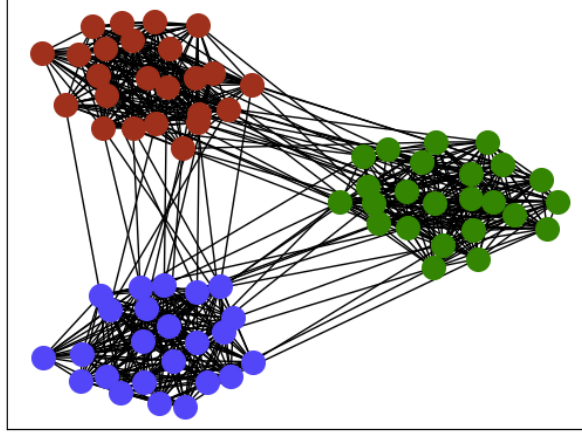


Figure 1: A realization of a Girvan-Newman benchmark graph.

Figure 2 shows the results of applying community detection algorithms on such graphs with varying values of the μ parameter. For each μ parameter value, the results were averaged over 25 iterations of graph construction and community detection. This method was used for all evaluations performed in this section.

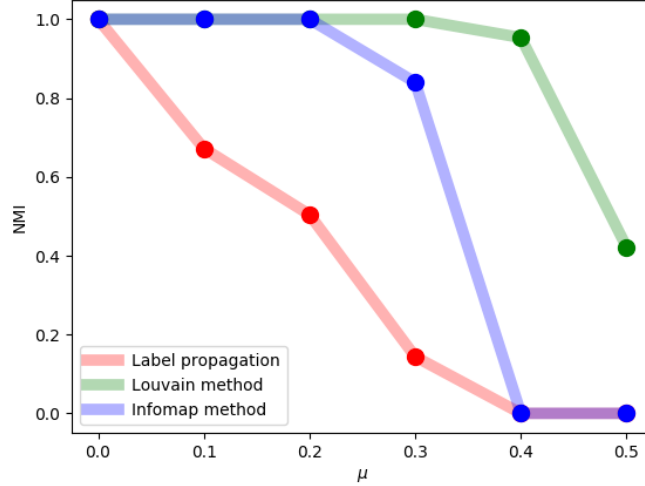


Figure 2: Normalized mutual information between the ground truth and results of community detection algorithms for varying values of the μ parameter of Girvan-Newman benchmark graphs.

We can see that the Louvain method performs best since it reliably detects the correct communities for values of μ lower than 0.4 and then drops sharply. Label propagation performs especially poorly as it decreases steadily for even moderately small values of the μ parameter. We can also see the clear absence of a desired sharp drop in normalized mutual information as the division into communities becomes too blurred to make reliable predictions. The Infomap method is stable up to moderate values of μ and then drops relatively sharply.

Figure 3 shows the results of performing such analysis on Lancichinetti synthetic benchmark graphs.

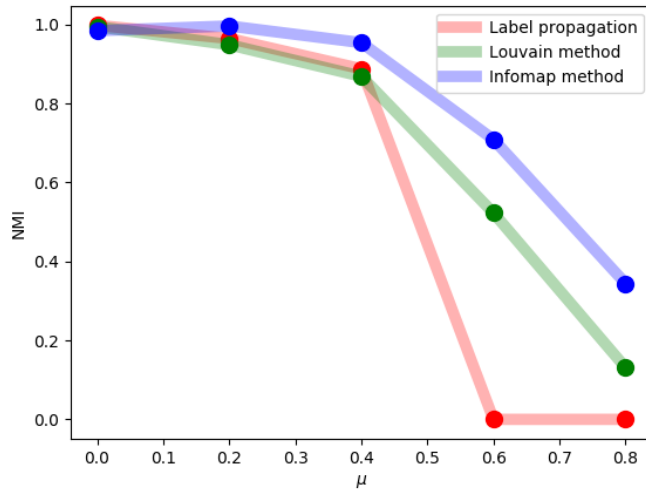


Figure 3: Normalized mutual information between the ground truth and results of community detection algorithms for varying values of the μ parameter of Lancichinetti synthetic benchmark graphs.

We can see that both the Louvain and Infomap method perform similarly and follow a very similar trajectory. All three methods are almost equivalent up to $\mu = 0.4$. The label propagation then drops sharply to 0 which can be seen as desirable at such large values of parameter μ . This makes the label propagation algorithm the algorithm with the most desirable properties when evaluated on such graphs.

Figure 4 shows the normalized variation of information computed using results of detecting communities in random Erdős–Rényi graph realizations, where each connected component corresponds to a community. The graphs are constructed with 1000 nodes and varying expected average degrees.

We can see that the label propagation method performs best as it best detects that the graph lacks communities and classifies the entire graph (its connected components) as single communities. The infomap method also performs similarly but in our case produces some false detections for very sparse random graphs. On the other hand, the Louvain method does not perform particularly well as it detects non-existent communities in the random graph realizations.

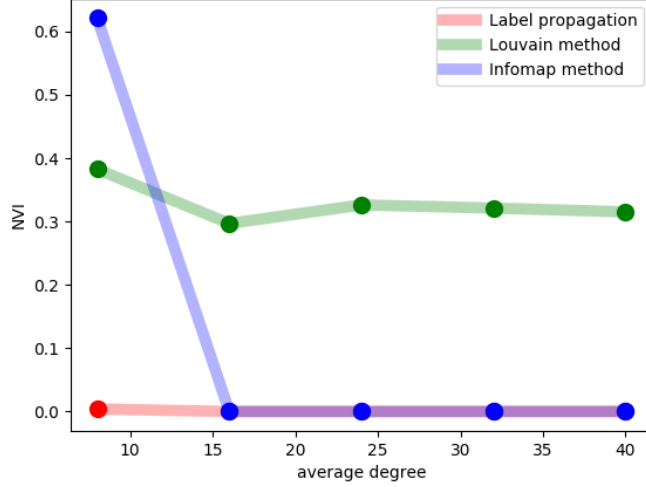


Figure 4: Normalized variation of informaton between the ground truth and results of community detection algorithms for varying expected node degree values of random Erdős–Rényi graph realizations.

Table 1 shows the normalized variation of information computed by considering pairs of runs of community detection algorithms on Lusseau bottlenose dolphins network. We can see that only the Louvain method produces results differing from zero. The other two algorithms perform deterministically and produce same results on each iteration.

| method | NVI |
|-------------------|--------|
| Label propagation | 0.0 |
| Louvain method | 0.1092 |
| Infomap method | 0.0 |

Table 1: Normalized variation of information of results obtained by pairs of sequential algorithm applications.

Overall, all three methods have their advantages and weaknesses. Their usefulness can be determined by studying their known properties and then choosing the method that is best considering what we know about our network. The label propagation method did not produce false detections in the Erdős–Rényi random graphs and had a desirable strong cut-off point when tested on Lancichinetti synthetic graphs. On the other hand, evaluations on simple Girvan-Newman graphs produced an undesirable almost linear drop in its performance that makes results interpretation very difficult. The Louvain method performed very well on Girvan-Newman and Lancichinetti benchmark graphs but suffered from false detections when used on Erdős–Rényi random graphs. In our experiments, the Infomap method performed relatively well on all benchmark graphs and the results indicate it to be the best all-round community detection method.

The code written to answer this question is given below.

```

1 import networkx as nx
2 import parse_network
3 import random
4 import matplotlib.pyplot as plt
5
6
7 def girvan_newman(num_groups, group_sizes, expected_degree, mu):
8     """
9     Construct Girvan-Newman benchmark graph with specified properties.
10
11     Args:
12         num_groups (int): Number of groups in the benchmark graph
13         group_sizes (int): Sizes of groups in the benchmark graph
14         expected_degree (int): expected node degree in the benchmark graph
15         mu (int): The mu parameter controlling the connectedness of the groups
16
17     Returns:
18         (tuple): Constructed graph and ground truth in required format
19     """
20
21     # Compute probability of a link between nodes in same group and
22     # link between nodes in different groups.
23     p_same = expected_degree*(1-mu)/(group_sizes-1)
24     p_other = expected_degree*mu/((num_groups-1)*group_sizes)
25
26     # Initialize empty graph with specified number of nodes.
27     graph = nx.empty_graph(n=num_groups*group_sizes, create_using=nx.Graph)
28
29     # Label groups of nodes and construct ground truth in required format.
30     nx.set_node_attributes(graph, {idx : {'label' : idx//group_sizes} for idx in range(graph
31         .number_of_nodes())})
32     attrs = nx.get_node_attributes(graph, 'label')
33     ground_truth = [{node_idx for node_idx, label in attrs.items() if label == comm_label}
34         for comm_label in set(attrs.values())]
35
36     # Add links.
37     node_idxs = list(graph.nodes())
38     for idx1 in range(len(node_idxs)-1):
39         for idx2 in range(idx1+1, len(node_idxs)):
40             if graph.nodes()[idx1]['label'] == graph.nodes()[idx2]['label']:
41                 # If nodes in same group, add link with probability p_same.
42                 if random.random() < p_same:
43                     graph.add_edge(idx1, idx2)
44             else:
45                 # If nodes not in same group, add link with probability p_other.
46                 if random.random() < p_other:
47                     graph.add_edge(idx1, idx2)
48
49     # Return constructed graph and ground truth.
50     return graph, ground_truth
51
52 def draw_girvan_newman(num_groups, group_sizes, expected_degree, mu):
53     """
54     Draw Girvan-Newman benchmark graph with specified properties and save plot to results
55     folder.
56
57     Args:
58         num_groups (int): Number of groups in the benchmark graph
59         group_sizes (int): Sizes of groups in the benchmark graph
60         expected_degree (int): expected node degree in the benchmark graph
61         mu (int): The mu parameter controlling the connectedness of the groups
62
63     Returns:
64         (int): 0 if success else 1.
65     """

```

```

65 # Initialize Girvan-Newman benchmark graph with specified properties.
66 graph, _ = girvan_newman(num_groups, group_sizes, expected_degree, mu)
67
68 # Get unique group labels.
69 labels = set(nx.get_node_attributes(graph, 'label').values())
70
71 # Set graph position (for plotting).
72 pos=nx.spring_layout(graph)
73
74 # Go over group labels and plot nodes corresponding to that group using random color.
75 for label in labels:
76     nodes_nxt_group = [n for (n, d) in graph.nodes(data=True) if d['label'] == label]
77     nx.draw_networkx_nodes(graph, pos, nodelist=nodes_nxt_group, node_size=200, node_
        color=[[random.random(), random.random(), random.random()]])
78
79 # Plot edges.
80 nx.draw_networkx_edges(graph, pos,width=1.0)
81
82 # Save figure.
83 try:
84     plt.savefig('../results/girvan_newman_benchmark_graph.png')
85     return 0
86 except:
87     return 1
88
89
90 def lancichinetti(mu):
91     """
92     Return Lancichinetti benchmark graph with specified mu parameter.
93
94     Args:
95         mu (float): The mu parameter.
96
97     Returns:
98         (tuple): Parsed graph with specified mu parameter and ground truth in required
99                 format
100     """
101     # Get path.
102     fmt = '{:<04}'
103     f_tmp = fmt.format(mu).replace('.', '')
104     f = 'LFR_' + f_tmp[:2] + '_' + f_tmp[2:]
105
106     # Load and parse graph. Get ground truth in required format.
107     graph = parse_network.parse_network('../data/LFR/' + f, create_using=nx.Graph)
108     attrs = nx.get_node_attributes(graph, 'data')
109     ground_truth = [{node_idx for node_idx, label in attrs.items() if label == comm_label}
        for comm_label in set(attrs.values())]
110
111     # Return graph and ground truth in required format.
112     return graph, ground_truth
113
114
115 def erdos_renyi(num_nodes, average_degree):
116     """
117     Construct Erdos-Renyi random graph with specified number of nodes and specified average
118         degree.
119
120     Args:
121         num_nodes (int): Number of nodes in constructed Erdos-Renyi random graph.
122         average_degree (int): Average degree in constructed Erdos-Renyi random graph
123
124     Returns:
125         (tuple): Parsed network and ground truth in required format
126     """
127     # Construct graph and return it along with its connected components as communities (
        ground truth).
    graph = nx.erdos_renyi_graph(num_nodes, average_degree/num_nodes, directed=False)

```



```

128     return graph, list(nx.algorithms.components.connected_components(graph))
129
130
131 def bottlenose_dolphins():
132     """
133     Parse and return Lusseau bottlenose dolphins network.
134
135     Returns:
136     (tuple): Parsed network and ground truth in required format
137     """
138
139     # Load and parse graph. Get ground truth in required format.
140     graph = parse_network.parse_network('../data/dolphins', create_using=nx.Graph)
141     attrs = nx.get_node_attributes(graph, 'data')
142     ground_truth = [{node_idx for node_idx, label in attrs.items() if label == comm_label}
143                     for comm_label in set(attrs.values())]
144
145     # Return graph and ground truth in required format.
146     return graph, ground_truth
147
148 if __name__ == '__main__':
149     # Draw Girvan-Newman benchmark graph and save plot.
150     draw_girvan_newman(3, 24, 20, 0.1)

```

```

1 import community
2 from cdlib import algorithms
3 import networkx as nx
4 import benchmark_graphs
5 import benchmark_utils
6 import pickle
7 import os
8 import sys
9
10
11 def benchmark_gn():
12     """
13     Perform benchmarking of algorithms on Girvan-Newman benchmark graph.
14
15     Returns:
16     (tuple): Used mu values, results for label propagation algorithm, results for
17             Louvain method,
18             results for Infomap method.
19     """
20
21     NUM_REP = 25 # number of algorithm repetitions (on newly constructed graph)
22     GN_NUM_GROUPS = 3 # number of groups in benchmark graph
23     GN_GROUP_SIZES = 24 # group sizes in benchmark graph
24     GN_EXPECTED_DEGREE = 20 # expected degree in benchmark graph
25     gn_mu_vals = (0.0, 0.1, 0.2, 0.3, 0.4, 0.5) # list of mu values for benchmark graph
26
27     # Initialize lists for storing results for different mu values.
28     y_vals_label_prop = []
29     y_vals_louvain = []
30     y_vals_infomap = []
31
32     # Go over mu values.
33     print("Performing benchmarks on Girvan-Newman benchmark graphs")
34     for idx, mu in enumerate(gn_mu_vals):
35
36         # Initialize lists for storing results for iterations.
37         nmi_label_prop = []
38         nmi_louvain = []
39         nmi_infomap = []
40
41         # Repeat benchmark graph construction and community detection specified number of

```

```

41         times.
42     for _ in range(NUM_REP):
43         # Construct benchmark graph with specified properties.
44         graph, ground_truth = benchmark_graphs.girvan_newman(GN_NUM_GROUPS, GN_GROUP_
45             SIZES, GN_EXPECTED_DEGREE, mu)
46
47         # Get detections for algorithms.
48         res_label_prop = benchmark_utils.normalize_community_format(nx.algorithms.
49             community.label_propagation.label_propagation_communities(graph), 'label_
50             propagation')
51         res_louvain = benchmark_utils.normalize_community_format(community.best_
52             partition(graph, randomize=True), 'louvain')
53         res_infomap = benchmark_utils.normalize_community_format(algorithms.infomap(
54             graph), 'infomap')
55
56         # Compute NMI values.
57         nmi_label_prop.append(benchmark_utils.nmi(res_label_prop, ground_truth))
58         nmi_louvain.append(benchmark_utils.nmi(res_louvain, ground_truth))
59         nmi_infomap.append(benchmark_utils.nmi(res_infomap, ground_truth))
60
61         # Get mean NMI value and set as value for current mu value.
62         y_vals_label_prop.append(sum(nmi_label_prop)/len(nmi_label_prop))
63         y_vals_louvain.append(sum(nmi_louvain)/len(nmi_louvain))
64         y_vals_infomap.append(sum(nmi_infomap)/len(nmi_infomap))
65         print("Done {0}/{1}".format(idx+1, len(gn_mu_vals)))
66
67     # Return data for plotting results.
68     return gn_mu_vals, y_vals_label_prop, y_vals_louvain, y_vals_infomap
69
70 def benchmark_lancichinetti():
71     """
72     Perform benchmarking of algorithms on Lancichinetti benchmark graph.
73
74     Returns:
75     (tuple): Used mu values, results for label propagation algorithm, results for
76         Louvain method,
77         results for Infomap method.
78     """
79
80     NUM_REP = 25 # number of algorithm repetitions (on newly constructed graph)
81     lanc_mu_vals = (0.0, 0.2, 0.4, 0.6, 0.8) # list of mu values for benchmark graph
82
83     # Initialize lists for storing results for different mu values.
84     y_vals_label_prop = []
85     y_vals_louvain = []
86     y_vals_infomap = []
87
88     # Go over mu values.
89     print("Performing benchmarks on Lancichinetti benchmark graphs")
90     for idx, mu in enumerate(lanc_mu_vals):
91         # Initialize lists for storing results for iterations.
92         nmi_label_prop = []
93         nmi_louvain = []
94         nmi_infomap = []
95
96         # Repeat benchmark graph construction and community detection specified number of
97         times.
98         for _ in range(NUM_REP):
99             # Construct benchmark graph with specified properties.
100             graph, ground_truth = benchmark_graphs.lancichinetti(mu)
101
102             # Get detections for algorithms.
103             res_label_prop = benchmark_utils.normalize_community_format(\
104                 nx.algorithms.community.label_propagation.label_propagation_communities(
105                     graph), 'label_propagation')

```

```

100     res_louvain = benchmark_utils.normalize_community_format(community.best_
101         partition(graph, randomize=True), 'louvain')
102     res_infomap = benchmark_utils.normalize_community_format(algorithms.infomap(
103         graph), 'infomap')
104
105     # Compute NMI values.
106     nmi_label_prop.append(benchmark_utils.nmi(res_label_prop, ground_truth))
107     nmi_louvain.append(benchmark_utils.nmi(res_louvain, ground_truth))
108     nmi_infomap.append(benchmark_utils.nmi(res_infomap, ground_truth))
109
110     # Get mean NMI value and set as value for current mu value.
111     y_vals_label_prop.append(sum(nmi_label_prop)/len(nmi_label_prop))
112     y_vals_louvain.append(sum(nmi_louvain)/len(nmi_louvain))
113     y_vals_infomap.append(sum(nmi_infomap)/len(nmi_infomap))
114     print("Done {0}/{1}".format(idx+1, len(lanc_mu_vals)))
115
116 # Return data for plotting results.
117 return lanc_mu_vals, y_vals_label_prop, y_vals_louvain, y_vals_infomap
118
119 def benchmark_er():
120     """
121     Perform benchmarking of algorithms on Erdos-Renyi random graph.
122
123     Returns:
124         (tuple): Used average node degrees, results for label propagation algorithm, results
125             for Louvain method,
126             results for Infomap method.
127     """
128
129     NUM_REP = 25 # number of algorithm repetitions (on newly constructed graph)
130     NUM_NODES = 1000 # number of nodes in benchmark graph
131     er_average_degrees = (8, 16, 24, 32, 40) # list of average degrees for benchmark graph
132
133     # Initialize lists for storing results for different mu values.
134     y_vals_label_prop = []
135     y_vals_louvain = []
136     y_vals_infomap = []
137
138     # Go over mu values.
139     print("Performing benchmarks on Erdos-Renyi random graphs")
140     for idx, av_deg in enumerate(er_average_degrees):
141
142         # Initialize lists for storing results for iterations.
143         nvi_label_prop = []
144         nvi_louvain = []
145         nvi_infomap = []
146
147         # Repeat benchmark graph construction and community detection specified number of
148         # times.
149         for _ in range(NUM_REP):
150
151             # Construct benchmark graph with specified properties.
152             graph, ground_truth = benchmark_graphs.erdos_renyi(NUM_NODES, av_deg)
153
154             # Get detections for algorithms.
155             res_label_prop = benchmark_utils.normalize_community_format(\
156                 nx.algorithms.community.label_propagation.label_propagation_communities(
157                     graph), 'label_propagation')
158             res_louvain = benchmark_utils.normalize_community_format(community.best_
159                 partition(graph, randomize=True), 'louvain')
160             res_infomap = benchmark_utils.normalize_community_format(algorithms.infomap(
161                 graph), 'infomap')
162
163             # Compute NMI values.
164             nvi_label_prop.append(benchmark_utils.nvi(res_label_prop, ground_truth))
165             nvi_louvain.append(benchmark_utils.nvi(res_louvain, ground_truth))
166             nvi_infomap.append(benchmark_utils.nvi(res_infomap, ground_truth))

```

```

161
162
163     # Get mean NMI value and set as value for current mu value.
164     y_vals_label_prop.append(sum(nvi_label_prop)/len(nvi_label_prop))
165     y_vals_louvain.append(sum(nvi_louvain)/len(nvi_louvain))
166     y_vals_infomap.append(sum(nvi_infomap)/len(nvi_infomap))
167     print("Done {0}/{1}".format(idx+1, len(er_average_degrees)))
168
169     # Return data for plotting results.
170     return er_average_degrees, y_vals_label_prop, y_vals_louvain, y_vals_infomap
171
172
173 def benchmark_dolphins():
174     """
175     Perform benchmarking of algorithms on Lusseau bottlenose dolphins network.
176
177     Returns:
178         (tuple): results for label propagation algorithm, results for Louvain method,
179         results for Infomap method.
180     """
181
182     NUM_REP = 25 # number of algorithm repetitions
183
184     # Initialize lists for storing results.
185     det_label_prop = []
186     det_louvain = []
187     det_infomap = []
188
189     # Repeat benchmark graph construction and community detection specified number of times.
190     print("Performing benchmarks on Lusseau bottlenose dolphins network")
191     for idx in range(NUM_REP):
192
193         # Construct benchmark graph with specified properties.
194         graph, ground_truth = benchmark_graphs.bottlenose_dolphins()
195
196         # Get detections for algorithms.
197         res_label_prop = benchmark_utils.normalize_community_format(\
198             nx.algorithms.community.label_propagation.label_propagation_communities(
199                 graph), 'label_propagation')
200         res_louvain = benchmark_utils.normalize_community_format(community.best_partition(
201             graph, randomize=True), 'louvain')
202         res_infomap = benchmark_utils.normalize_community_format(algorithms.infomap(graph),
203             'infomap')
204
205         # Add detections to results list.
206         det_label_prop.append(res_label_prop)
207         det_louvain.append(res_louvain)
208         det_infomap.append(res_infomap)
209
210         print("Done {0}/{1}".format(idx+1, NUM_REP))
211
212     # Initialize lists for computing pairwise NVI for detections.
213     pairwise_nvi_label_prop = []
214     pairwise_nvi_louvain = []
215     pairwise_nvi_infomap = []
216
217     # Compute NVI of detections (pairwise).
218     for idx in range(NUM_REP-1):
219         pairwise_nvi_label_prop.append(benchmark_utils.nvi(det_label_prop[idx], det_label_
220             prop[idx+1]))
221         pairwise_nvi_louvain.append(benchmark_utils.nvi(det_louvain[idx], det_louvain[idx
222             +1]))
223         pairwise_nvi_infomap.append(benchmark_utils.nvi(det_infomap[idx], det_infomap[idx
224             +1]))
225
226     # Get mean NVI values and set as results
227     res_label_prop = sum(pairwise_nvi_label_prop)/len(pairwise_nvi_label_prop)
228     res_louvain = sum(pairwise_nvi_louvain)/len(pairwise_nvi_louvain)

```

```

223     res_infomap = sum(pairwise_nvi_infomap)/len(pairwise_nvi_infomap)
224
225     # Return data for plotting results.
226     return res_label_prop, res_louvain, res_infomap
227
228
229 def plot_results(x, y, labels_y, x_label, y_label, file_name):
230     """
231     Plot benchmarking results.
232
233     Args:
234         x (list): x-axis values
235         y (list): y-axis values (list of lists)
236         labels_y (list): Labels for the drawn lines (for legend)
237         x_label (str): The x-axis label
238         y_label (str): The y-axis label
239         file_name (str): The file name for the saved plot
240     """
241
242     import matplotlib.pyplot as plt
243
244     # Plot results.
245     fig = plt.figure()
246     plt.plot(x, y[0], 'r.', markersize=20)
247     plt.plot(x, y[0], 'r-', label=labels_y[0], alpha=0.3, linewidth=7)
248     plt.plot(x, y[1], 'g.', markersize=20)
249     plt.plot(x, y[1], 'g-', label=labels_y[1], alpha=0.3, linewidth=7)
250     plt.plot(x, y[2], 'b.', markersize=20)
251     plt.plot(x, y[2], 'b-', label=labels_y[2], alpha=0.3, linewidth=7)
252     plt.xlabel(x_label)
253     plt.ylabel(y_label)
254     plt.legend()
255
256     # Save plot to file.
257     plt.savefig('../results/' + file_name)
258
259
260 def perform_benchmarking():
261     """
262     Perform implemented benchmarks.
263
264     Returns:
265         (int): 0 if successful else 1
266     """
267
268     try:
269
270         # Compute results using benchmarking functions.
271         gn_mu_vals, y_vals_label_prop_gn, y_vals_louvain_gn, y_vals_infomap_gn = benchmark_
            gn()
272         lanc_mu_vals, y_vals_label_prop_lc, y_vals_louvain_lc, y_vals_infomap_lc = benchmark
            _lancichinetti()
273         er_average_degrees, y_vals_label_prop_er, y_vals_louvain_er, y_vals_infomap_er =
            benchmark_er()
274         res_label_prop_dolph, res_louvain_dolph, res_infomap_dolph = benchmark_dolphins()
275
276         # Build dictionary for results.
277         res = {'gn_mu_vals' : gn_mu_vals,
278               'y_vals_label_prop_gn' : y_vals_label_prop_gn,
279               'y_vals_louvain_gn' : y_vals_louvain_gn,
280               'y_vals_infomap_gn' : y_vals_infomap_gn,
281               'lanc_mu_vals' : lanc_mu_vals,
282               'y_vals_label_prop_lc' : y_vals_label_prop_lc,
283               'y_vals_louvain_lc' : y_vals_louvain_lc,
284               'y_vals_infomap_lc' : y_vals_infomap_lc,
285               'er_average_degrees' : er_average_degrees,
286               'y_vals_label_prop_er' : y_vals_label_prop_er,
287               'y_vals_louvain_er' : y_vals_louvain_er,

```

```

288         'y_vals_infomap_er' : y_vals_infomap_er,
289         'res_label_prop_dolph' : res_label_prop_dolph,
290         'res_louvain_dolph' : res_louvain_dolph,
291         'res_infomap_dolph' : res_infomap_dolph,
292     }
293
294     # Save dictionary.
295     with open('../results/cached_data_3.p', 'wb') as f:
296         pickle.dump(res, f, pickle.HIGHEST_PROTOCOL)
297     return 0
298 except:
299     return 1
300
301 if __name__ == '__main__':
302     import argparse
303
304     # Parse optional flag for parsing cached results.
305     parser = argparse.ArgumentParser()
306     parser.add_argument('--load-cached', action='store_true')
307     args = parser.parse_args()
308
309     # If not loading cached results, perform benchmarking.
310     if not args.load_cached:
311         res = perform_benchmarking()
312         if res == 0:
313             print("Benchmarking completed")
314         else:
315             print("Something went wrong during the benchmarking process")
316
317     if os.path.isfile('../results/cached_data_3.p'):
318
319         # If cached data exists, load it and visualize results.
320         with open('../results/cached_data_3.p', 'rb') as f:
321             res = pickle.load(f)
322
323         # Plot results for Girvan-Newman benchmark graph.
324         plot_results(res['gn_mu_vals'], [res['y_vals_label_prop_gn'], res['y_vals_louvain_gn'],
325             res['y_vals_infomap_gn']],
326             ['Label propagation', 'Louvain method', 'Infomap method'], x_label=r'$\mu$',
327             y_label='NMI', file_name='benchmark_gn.png')
328
329         # Plot results for Lancichinetti benchmark graph.
330         plot_results(res['lanc_mu_vals'], [res['y_vals_label_prop_lc'], res['y_vals_louvain_lc'],
331             res['y_vals_infomap_lc']],
332             ['Label propagation', 'Louvain method', 'Infomap method'], x_label=r'$\mu$',
333             y_label='NMI', file_name='benchmark_lc.png')
334
335         # Plot results for Erdos-Renyi benchmark (random) graph.
336         plot_results(res['er_average_degrees'], [res['y_vals_label_prop_er'], res['y_vals_louvain_er'],
337             res['y_vals_infomap_er']],
338             ['Label propagation', 'Louvain method', 'Infomap method'], x_label='average degree',
339             y_label='NVI', file_name='benchmark_er.png')
340
341         # Save results for dolphin network in form of Markdown table.
342         with open('../results/res_dolphins.txt', 'w') as f:
343             f.write('| method | NVI | \n')
344             f.write('|-----|-----| \n')
345             f.write('| Label propagation | {0:.4f} | \n'.format(res['res_label_prop_dolph']))
346             f.write('| Louvain method | {0:.4f} | \n'.format(res['res_louvain_dolph']))
347             f.write('| Infomap method | {0:.4f} | \n'.format(res['res_infomap_dolph']))
348
349     sys.exit(0)
350 else:
351     print("Cached data does not exist. Please run script without the --load-cached flag.")
352     sys.exit(1)

```

```

1 import math
2 import scipy as sp
3 from sklearn.metrics import normalized_mutual_info_score
4 from pyitlib import discrete_random_variable as drv
5
6
7 def nmi(prediction, ground_truth):
8     """
9     Compute normalized mutual information of predicted communities with ground truth.
10
11     Args:
12         prediction (list): List of sets representing found communities
13         ground_truth (list): List of sets representing ground truth
14
15     Returns:
16         (float): Computed normalized mutual information score
17     """
18
19     # Assign labels to nodes, sort by node index and get labels.
20     first_partition_c = [x[1] for x in sorted([(node, nid) for nid, cluster in enumerate(
21         prediction) for node in cluster], key=lambda x: x[0])]
22     second_partition_c = [x[1] for x in sorted([(node, nid) for nid, cluster in enumerate(
23         ground_truth) for node in cluster], key=lambda x: x[0])]
24
25     # Compute normalized mutual information.
26     return normalized_mutual_info_score(first_partition_c, second_partition_c)
27
28 def nvi(prediction, ground_truth):
29     """
30     Compute normalized variation of information of predicted communities with ground truth.
31
32     Args:
33         prediction (list): List of sets representing found communities
34         ground_truth (list): List of sets representing ground truth
35
36     Returns:
37         (float): Computed normalized variation of information score
38     """
39
40     # Assign labels to nodes, sort by node index and get labels.
41     first_partition_c = [x[1] for x in sorted([(node, nid) for nid, cluster in enumerate(
42         prediction) for node in cluster], key=lambda x: x[0])]
43     second_partition_c = [x[1] for x in sorted([(node, nid) for nid, cluster in enumerate(
44         ground_truth) for node in cluster], key=lambda x: x[0])]
45
46     # Compute normalized mutual information.
47     return drv.information_variation(first_partition_c, second_partition_c, base=math.e)/
48         math.log(len(first_partition_c))
49
50
51 def normalize_community_format(res, method):
52     """
53     Normalize predicted communities format for all methods to list of sets.
54
55     Args:
56         res (list): Predicted communities as returned by method
57         method (str): Method used to predict the communities ('label_propagation', 'louvain', 'infomap')
58
59     Returns:
60         (list): List of sets representing found communities
61     """
62
63     if method == 'label_propagation':
64         return list(res)
65     elif method == 'louvain':

```

```

62         return [{node_idx for node_idx, label in res.items() if label == com_label} for com_
        label in set(res.values())]
63     elif method == 'infomap':
64         return list(map(lambda x: set(x), res.communities))

```

```

1 import networkx as nx
2 import re
3
4
5 def parse_network(path, *args, **kwargs):
6     """
7     Parse network and add associated data. The data should be specified using the LNA format
8     .
9     Author: Jernej Vivod
10
11     Args:
12         path (str): Path to the data file.
13         args (tuple): Arguments for the networkx read_edgelist method.
14         **kwargs (dict): Keyword arguments for the networkx read_edgelist method.
15
16     Returns:
17         (obj): Networkx graph representation with added node names and data.
18     """
19
20     def parse_line_data(line):
21         node_idx = line[:line.index(',')].split(' ')[1]
22         node_name = re.findall(r'"(["]*)"', line)[0]
23         node_data = line.split('" "')[1].strip()
24         return node_idx, node_name, node_data
25
26     # Parse graph from edge list.
27     graph = nx.read_edgelist(path, *args, **kwargs)
28
29     # Initialize dictionaries for parsing data.
30     names = dict()
31     data = dict()
32
33     # Flag indicating the start of parsing.
34     parse = False
35     with open(path, 'r') as f:
36
37         # Go over lines.
38         for line in f:
39
40             # If parsing.
41             if parse:
42                 if len(line.split(" ")) == 1:
43                     # If found last delimiter, add data to graph and return.
44                     nx.set_node_attributes(graph, names, 'name')
45                     nx.set_node_attributes(graph, data, 'data')
46                     return graph
47                 else:
48                     # Parse node index, name and associated data and add to dictionaries.
49                     node_idx, node_name, node_data = parse_line_data(line)
50                     data[node_idx] = node_data
51                     names[node_idx] = node_name
52             else:
53                 if len(line.split(" ")) == 1:
54                     # If at first delimiter, set parse flag.
55                     parse = True
56                 else:
57                     pass

```


Question 4

Answer the given question concerning link prediction in real-world networks. Implement a framework for evaluating link prediction methods and use it in conjunction with specified sample networks.

We know that real networks are usually sparse. A single node is only connected to a few other nodes. Real networks are also usually large. Assuming these two properties, a link between two randomly chosen nodes in such a network is very unlikely. The classification accuracy of a method which simply predicts that no link will occur will therefore tend to 1.0 as the size of the network increases.

Table 2 shows the average AUC obtained by running the algorithms on each network 10 times.

| method | preferential attachment | Adamic-Adar | Community |
|-------------|-------------------------|-------------|-----------|
| Erdős-Rényi | 0.5109 | 0.5006 | 0.6548 |
| Gnutella | 0.7169 | 0.5128 | 0.7866 |
| Facebook | 0.8292 | 0.9925 | 0.9513 |
| nec | 0.8181 | 0.6854 | 0.9280 |

Table 2: AUC scores obtained by benchmarking different link prediction methods on sample networks.

For the Erdős-Rényi graph, the community algorithm proves most successful. As all connected components of the random graph consist of a single community, the indices will be the same for each node as long as the community detection algorithm successfully classifies the connected components as single communities. This can be problematic when using the Louvain method. The equal indices best capture the random structure of this kind of graph. For the Gnutella network, the community index and the preferential attachment models perform well. By plotting a portion of the network, we can clearly see, that nodes tend to form communities centered around hubs that have a few edges leading out to other communities. The degree distribution of the network also follows the power-law which is why the preferential attachment assumption also yields satisfactory results. The absence of triadic closure makes the Adamic-Adar index unsuitable for such a network. On the other hand, the Adamic-Adar index performs very well on the Facebook social circles networks as triadic closure or abundance of triangles is a notable property of such social networks. The network also follows a power-law degree distribution and contains notable communities, which is why the preferential attachment index and the community index also perform relatively well. Similarly as the Gnutella network, the nec overlay map lacks triadic closures, which is to be expected for these types of networks as such connections would introduce redundancy. The degree distribution again follows an approximate power-law distribution and the nodes again form communities with central hubs which is why the preferential attachment and community indices perform relatively well with the community index best capturing the nature of this particular network.

The code written to answer this question is given below.

```

1 import random
2 import math
3 import community
4 from scipy.special import comb
5 from collections import Counter
6 import parse_network
7
8
9 def link_prediction_auc(network, prediction_func):
10     """
11     Perform link prediction using specified method on specified network and
12     return AUC.
13
14     Args:
15         network (object): The network on which to evaluate the link prediction mechanism
16         prediction_func (function): The function implementing link prediction index
17         computation
18
19     Returns:
20         (float): AUC score of method
21     """
22     # Randomly sample m/10 pairs of nodes that are not yet
23     # linked and store them into L_{N}.
24     negative_examples = []
25     while len(negative_examples) < math.ceil(network.number_of_edges()/10):
26         pair = tuple(random.sample(network.nodes, 2))
27         if not network.has_edge(*pair):
28             negative_examples.append(pair)
29
30     # Randomly sample m/10 links from the network, remove them from the
31     # network and store them into L_{P}.
32     positive_examples = random.sample(network.edges, math.ceil(network.number_of_edges()/10))
33     network.remove_edges_from(positive_examples)
34
35     # Compute the link prediction index s for all pairs of nodes in union of L_{N}
36     # and L_{P}.
37     # link_prediction_indices = lp_idx(negative_examples + positive_examples)
38     lp_ind_n = [prediction_func(network, link) for link in negative_examples]
39     lp_ind_p = [prediction_func(network, link) for link in positive_examples]
40
41     # Sample m/10 pairs from L_{N} and m/10 pairs from L_{P} with repetitions.
42     samp_lp_ind_n = random.choices(lp_ind_n, k=int(math.ceil(network.number_of_edges()/10)))
43     samp_lp_ind_p = random.choices(lp_ind_p, k=int(math.ceil(network.number_of_edges()/10)))
44     comp = [val_p - val_n for val_p, val_n in zip(samp_lp_ind_p, samp_lp_ind_n)]
45     num_p_larger = sum(e1 > 0 for e1 in comp)
46     num_eq = sum(e1 == 0 for e1 in comp)
47
48     return (num_p_larger + num_eq/2)/(math.ceil(network.number_of_edges()/10))
49
50
51 def get_index_func(kind, network):
52     """
53     Get specified function for computing link prediction index.
54
55     Args:
56         kind (str): Which link prediction index function to return
57         network (object): The network for which link prediction will be run
58
59     Returns:
60         (function): Specified function for computing link prediction index.
61     """
62
63     ### Function used to compute the link indices ###
64
65

```

```

66 # Compute preferential attachment index.
67 def preferential_attachment_index(network, link):
68     return network.degree[link[0]]*network.degree[link[1]]
69
70 # Compute Adamic-Adar index.
71 def adamic_adar_index(network, link):
72     return sum(1/math.log(network.degree(x))
73               for x in set(network.neighbors(link[0])).intersection(network.neighbors(link
74                               [1])))
75
76 # Compute community index.
77 def community_index(network, communities, nc, mc, link):
78     if communities[link[0]] != communities[link[1]]:
79         return 0
80     else:
81         return mc[communities[link[0]]]/comb(nc[communities[link[0]]], 2)
82
83 #####
84
85 # Get counts of edges in communities.
86 def get_mc(network, communities):
87     counts = dict.fromkeys(set(communities.values()), 0)
88     for edge in network.edges():
89         if communities[edge[0]] == communities[edge[1]]:
90             counts[communities[edge[0]]] += 1
91     return counts
92
93 # Return specified function.
94 if kind == 'preferential-attachment':
95     return preferential_attachment_index
96 elif kind == 'adamic-adar':
97     return adamic_adar_index
98 elif kind == 'community':
99     communities = community.best_partition(network)
100     nc = Counter(communities.values())
101     mc = get_mc(network, communities)
102     return lambda network, link: community_index(network, communities, nc, mc, link)
103 else:
104     raise(ValueError('Unknown index function specified.'))
105
106
107 def get_benchmark_network(name):
108     """
109     Get benchmark networks for comparing link prediction methods.
110
111     Args:
112         name (str): Name of the benchmark network to get
113
114     Returns:
115         (object): Parsed network
116     """
117
118     if name == 'erdos-renyi':
119         # Erdos-Renyi random graph
120         NUM_NODES = 25000
121         AVERAGE_DEGREE = 10
122         from benchmark_graphs import erdos_renyi
123         network, _ = erdos_renyi(num_nodes=NUM_NODES, average_degree=AVERAGE_DEGREE)
124         return network
125
126     elif name == 'gnutella':
127         # Gnutella peer-to-peer file sharing network
128         return parse_network.parse_network('../data/gnutella', create_using=nx.Graph)
129
130     elif name == 'facebook':
131         # Facebook social circles network
132         return parse_network.parse_network('../data/circles', create_using=nx.Graph)

```

```

133
134     elif name == 'nec':
135         # nec overlay map
136         return parse_network.parse_network('../data/nec', create_using=nx.Graph)
137
138     else:
139         raise(ValueError('Unknown network specified'))
140
141
142 def main():
143     """
144     Perform benchmarking of link prediction methods and save results.
145     """
146
147     # Number of runs of each method to perform.
148     NUM_RUNS = 1
149
150     ### RUN EVALUTATIONS AND WRITE RESULTS TO FILE ###
151     # Initialize lists for storing results for runs. Perform evaluations and save average
152     # result.
153
154     auc_pref_vals = []
155     auc_ad_vals = []
156     auc_comm_vals = []
157     network = get_benchmark_network('erdos-renyi')
158     for idx in range(NUM_RUNS):
159         auc_pref_vals.append(link_prediction_auc(network, get_index_func('preferential-
160             attachment', network)))
161         auc_ad_vals.append(link_prediction_auc(network, get_index_func('adamic-adar',
162             network)))
163         auc_comm_vals.append(link_prediction_auc(network, get_index_func('community',
164             network)))
165     with open('../results/res_auc.txt', 'w') as f:
166         f.write("|          | preferential attachment | adamic-Adar | community |\n")
167         f.write("|-----|-----|-----|-----|\n")
168         f.write("| Erdos-Renyi | {0:.4f}          | {1:.4f}          | {2:.4f}          |\n".
169             format(sum(auc_pref_vals)/len(auc_pref_vals),
170                 sum(auc_ad_vals)/len(auc_ad_vals), sum(auc_comm_vals)/len(auc_comm_vals)))
171
172     auc_pref_vals = []
173     auc_ad_vals = []
174     auc_comm_vals = []
175     network = get_benchmark_network('gnutella')
176     for idx in range(NUM_RUNS):
177         auc_pref_vals.append(link_prediction_auc(network, get_index_func('preferential-
178             attachment', network)))
179         auc_ad_vals.append(link_prediction_auc(network, get_index_func('adamic-adar',
180             network)))
181         auc_comm_vals.append(link_prediction_auc(network, get_index_func('community',
182             network)))
183     with open('../results/res_auc.txt', 'a') as f:
184         f.write("| gnutella          | {0:.4f}          | {1:.4f}          | {2:.4f}          |\n".
185             format(sum(auc_pref_vals)/len(auc_pref_vals),
186                 sum(auc_ad_vals)/len(auc_ad_vals), sum(auc_comm_vals)/len(auc_comm_vals)))
187
188     auc_pref_vals = []
189     auc_ad_vals = []
190     auc_comm_vals = []
191     network = get_benchmark_network('facebook')
192     for idx in range(NUM_RUNS):
193         auc_pref_vals.append(link_prediction_auc(network, get_index_func('preferential-
194             attachment', network)))
195         auc_ad_vals.append(link_prediction_auc(network, get_index_func('adamic-adar',
196             network)))
197         auc_comm_vals.append(link_prediction_auc(network, get_index_func('community',
198             network)))

```

```

189 with open('../results/res_auc.txt', 'a') as f:
190     f.write("| facebook | {0:.4f} | {1:.4f} | {2:.4f} | \n".
191             format(sum(auc_pref_vals)/len(auc_pref_vals),
192                     sum(auc_ad_vals)/len(auc_ad_vals), sum(auc_comm_vals)/len(auc_comm_vals)))
193
194 auc_pref_vals = []
195 auc_ad_vals = []
196 auc_comm_vals = []
197 network = get_benchmark_network('nec')
198 for idx in range(NUM_RUNS):
199     auc_pref_vals.append(link_prediction_auc(network, get_index_func('preferential-
200     attachment', network)))
201     auc_ad_vals.append(link_prediction_auc(network, get_index_func('adamic-adar',
202     network)))
203     auc_comm_vals.append(link_prediction_auc(network, get_index_func('community',
204     network)))
205
206 with open('../results/res_auc.txt', 'a') as f:
207     f.write("| nec | {0:.4f} | {1:.4f} | {2:.4f} | \n".
208             format(sum(auc_pref_vals)/len(auc_pref_vals),
209                     sum(auc_ad_vals)/len(auc_ad_vals), sum(auc_comm_vals)/len(auc_comm_vals)))
210
211 if __name__ == '__main__':
212     main()

```

Question 5

Propose a strategy of predicting journals corresponding to physics papers that gives at least $\approx 70\%$ classification accuracy. The strategy can use any network analysis method or other approach as long as it scales better than $\mathcal{O}(n^2)$ in real networks.

We can solve this task by splitting the dataset (our network) into training and test sets. In this case, the training set consists of nodes corresponding to papers published in years before 2013 and the test set of papers published in year 2013. We can extract useful features from these nodes such as the node degree, the mean degree of its neighbours, the maximum degree of its neighbours, the minimum degree of its neighbours, the standard deviation of the degree of its neighbours, the number of neighbours with the same label, the number of triangles the node is part of, the maximum number of triangles a neighbour is part of, the minimum number of triangles the neighbour is part of, the mean number of triangles a neighbour is part of and the number of times each label (journal title) appears in the neighbourhood.

This yields a vector describing each node in the network. We can compute such features for each node in the training and test sets and get a canonical matrix form where each row corresponds to a node.

Standard machine learning algorithms can now be used to predict the node labels (titles of journals). We used a random-forest classifier to perform the classifications. We compared the results with a baseline classifier that predicts the label (corresponding journal) of each node to be the most common label in its neighbourhood.

The average classification accuracy of our approach is **0.88**. The average classification accuracy of the baseline approach is **0.66**.

Table 3 shows the classification accuracies and other statistics of our approach. Table 4 shows the same statistics for the baseline approach. We can clearly see that our approach significantly outperforms the baseline.

| Label | Precision | Recall | f1-score | Support |
|-------------|-----------|--------|----------|---------|
| PhysRevA | 0.86 | 0.89 | 0.88 | 2385 |
| PhysRevB | 0.87 | 0.95 | 0.91 | 4090 |
| PhysRevC | 0.91 | 0.93 | 0.92 | 883 |
| PhysRevD | 0.96 | 0.96 | 0.96 | 2868 |
| PhysRevE | 0.86 | 0.80 | 0.83 | 1875 |
| PhysRevLett | 0.86 | 0.80 | 0.83 | 3336 |
| PhysRevSTAB | 0.89 | 0.74 | 0.81 | 129 |
| PhysRevX | 0.81 | 0.26 | 0.39 | 86 |
| RevModPhys | 0.96 | 0.51 | 0.67 | 43 |

Table 3: Precision, recall, f1-score and support for our approach.

| Label | Precision | Recall | f1-score | Support |
|-------------|-----------|--------|----------|----------|
| PhysRevA | 0.77 | 0.61 | 0.880.68 | 23852385 |
| PhysRevB | 0.79 | 0.62 | 0.69 | 4090 |
| PhysRevC | 0.86 | 0.82 | 0.84 | 883 |
| PhysRevD | 0.88 | 0.89 | 0.89 | 2868 |
| PhysRevE | 0.85 | 0.52 | 0.64 | 1875 |
| PhysRevLett | 0.36 | 0.60 | 0.45 | 3336 |
| PhysRevSTAB | 0.80 | 0.60 | 0.69 | 129 |
| PhysRevX | 0.00 | 0.00 | 0.00 | 86 |
| RevModPhys | 0.00 | 0.00 | 0.00 | 43 |

Table 4: Precision, recall, f1-score and support for baseline approach.

The code written to answer this question is given below.

```

1
2 import networkx as nx
3 import numpy as np
4 import re
5 import parse_network
6 from sklearn import preprocessing
7 from sklearn.metrics import classification_report
8
9
10 def get_tts(network):
11     """
12     Get train-test split for network nodes corresponding to papers published in 2013
13     are added to the test set and nodes corresponding to papers published earlier are
14     added to the training set.
15
16     Args:
17         network (object): Network on which to perform the train-test split.
18
19     Returns:

```

```

20     (tuple): Indices of nodes training nodes and indices of test nodes.
21     """
22
23     # Find nodes corresponding to papers published in 2013 (test set) and nodes
24     # corresponding to papers
25     # published earlier (train set).
26     train_idx = [idx for idx, data in dict(network.nodes(data=True)).items() if data['name']
27                  ][-4:] != '2013']
28     test_idx = [idx for idx, data in dict(network.nodes(data=True)).items() if data['name']
29                 ][-4:] == '2013']
30
31     # Return training data indices and test data indices.
32     return train_idx, test_idx
33
34 def get_features_node(node, network, bow):
35     """
36     Get features for specified node.
37
38     Args:
39         node (str): Node index
40         network (object): The network the node is part of
41         bow (list): List of all node labels
42
43     Returns:
44         (numpy.ndarray): Vector of features for the current node
45     """
46
47     # Get paper of node.
48     name = network.node[node]['name']
49
50     # Get target value.
51     target = re.findall('[a-zA-Z]+', name)[0]
52
53     # Compute degree, mean degree of neighbors, number of triangles the
54     # node forms with neighbors.
55     degree = network.degree[node]
56     neighbors = [n for n in network.neighbors(node)]
57     neighbors_names = [re.findall('[a-zA-Z]+', network.node[neigh]['name'])[0] for neigh in
58                       neighbors]
59
60     # Compute bag-of-words features (number of neighbors with each label).
61     bow_feature = np.zeros(len(bow), dtype=int)
62     for idx, w in enumerate(bow):
63         bow_feature[idx] = neighbors_names.count(w)
64
65     # Compute mean degree of neighbors.
66     mean_degree_neigh = np.mean([network.degree[neigh] for neigh in neighbors])
67
68     # Compute maximum degree of neighbors.
69     max_degree_neigh = np.max([network.degree[neigh] for neigh in neighbors])
70
71     # Compute minimum degree of neighbors.
72     min_degree_neigh = np.min([network.degree[neigh] for neigh in neighbors])
73
74     # Compute standard deviation of degree of neighbors.
75     std_degree_neigh = np.std([network.degree[neigh] for neigh in neighbors])
76
77     # Compute number of neighbors with same target.
78     num_neighbors_same_target = 0
79     for neigh in neighbors:
80         if re.findall('[a-zA-Z]+', network.node[neigh]['name'])[0] == target:
81             num_neighbors_same_target += 1
82
83     # Compute number of triangles including current nodes.
84     num_triangles_this = nx.triangles(network, node)

```

```

84     # Compute mean number of triangles including one of the neighbors.
85     triangles_neighbors = [nx.triangles(network, neigh) for neigh in neighbors]
86
87     # Compute mean number of triangles of neighbors.
88     mean_triangles_neigh = np.mean(triangles_neighbors)
89
90     # Compute maximum number of triangles of neighbors.
91     max_triangles_neigh = np.max(triangles_neighbors)
92
93     # Compute minimum number of triangles of neighbors.
94     min_triangles_neigh = np.min(triangles_neighbors)
95
96     # Construct features vectors.
97     feature_vec = np.append(np.array([degree,
98                                     mean_degree_neigh,
99                                     max_degree_neigh,
100                                    min_degree_neigh,
101                                    std_degree_neigh,
102                                    num_neighbors_same_target,
103                                    num_triangles_this,
104                                    mean_triangles_neigh,
105                                    max_triangles_neigh,
106                                    min_triangles_neigh]), bow_feature)
107
108     # Return features vector and target variable.
109     return feature_vec, target
110
111
112 def get_features(network, node_idxs):
113     """
114     Get features for specified node indices.
115
116     Args:
117         network (object): The network containing the nodes
118         node_idxs (list): List of node indices for which to compute
119                             features
120
121     Returns:
122         (tuple): numpy array containing features and numpy array
123                 containing the target variable values.
124     """
125
126     # Define and initialize data and target variables.
127     data = None
128     target = []
129
130     # Get label encoder and "bag-of-words".
131     le, bow = get_label_encoder_and_bow(network)
132
133     # Go over specified nodes and compute features.
134     for idx, node in enumerate(node_idxs):
135         print('done {0}/{1}'.format(idx, len(node_idxs)))
136         feature_vec_nxt, target_nxt = get_features_node(node, network, bow)
137         target.append(target_nxt)
138         if data is None:
139             data = feature_vec_nxt
140         else:
141             data = np.vstack((data, feature_vec_nxt))
142
143     # Perform label encoding for target variable.
144     target = le.transform(target)
145
146     # Return data and target arrays.
147     return data, target
148
149
150 def get_label_encoder_and_bow(network):
151     """

```



```

152     Get label encoder for target variables and "bag-of-words".
153
154     Args:
155         (network): network from which to take the labels.
156
157     Returns:
158         (tuple): Fitted LabelEncoder instance and "bag-of-words" list sorted in
159                 alphabetical order
160     """
161
162     # Get labels found in network.
163     names = nx.get_node_attributes(network, 'name').values()
164
165     # Get unique labels as "bag-of-words".
166     bow = list(map(lambda x: re.findall('[a-zA-Z]+', x)[0], names))
167
168     # Fit label-encoder on bag-of-words.
169     le = preprocessing.LabelEncoder().fit(bow)
170
171     # Return fitted label encoder and sorted "bag-of-words".
172     return le, sorted(list(set(bow)))
173
174
175 def majority_neigh(network, node_idx):
176     """
177     Get results for baseline classifier that predicts the label to be the most
178     common label among the neighbors.
179
180     Args:
181         network (object): The network containing the nodes
182         node_idx (list): List of node indices for which to make the predictions
183
184     Returns:
185         (list): List of label predictions for nodes in node_idx list
186     """
187
188     # Get label encoder.
189     le, _ = get_label_encoder_and_bow(network)
190
191     # Initialize list for storing the results.
192     res = []
193
194     # Go over nodes and perform classification.
195     for node in node_idx:
196         neighbors = network.neighbors(node)
197         neighbor_targets = [re.findall('[a-zA-Z]+', network.node[neigh]['name'])[0] for
198                             neigh in neighbors]
199         res.append(max(set(neighbor_targets), key=neighbor_targets.count))
200
201     # Return classifications.
202     return le.transform(res)
203
204 def evaluate_model(data_train, target_train, data_test, target_test, clf):
205     """
206     Evaluate model using specified classifier.
207
208     Args:
209         data_train (numpy.ndarray): Training data
210         target_train (numpy.ndarray): Training target values
211         data_test (numpy.ndarray): Test data
212         target_test (numpy.ndarray): Test training variables
213         clf (object): Classifier to evaluate
214
215     Returns:
216         (str): Classification report
217     """
218

```

```

219     # Fit classifier.
220     clf.fit(data_train, target_train)
221
222     # Score predictions and create classification report.
223     pred = clf.predict(data_test)
224     return classification_report(target_test, pred)
225
226
227 def main():
228     """
229     Split data into training and test sets, construct features and evaluate classifiers.
230     """
231
232     from sklearn.ensemble import RandomForestClassifier
233
234
235     # Parse network.
236     network = parse_network.parse_network('../data/aps_2008_2013', create_using=nx.Graph)
237
238     import pdb
239     pdb.set_trace()
240
241     le, _ = get_label_encoder_and_bow(network)
242
243     # Get training and test data.
244     train_idx, test_idx = get_tts(network)
245     data_train, target_train = get_features(network, train_idx)
246     data_test, target_test = get_features(network, test_idx)
247
248     # Initialize random-forest classifier.
249     clf_rf = RandomForestClassifier()
250
251     # Compute classification reports and write to file.
252     clf_report_rf = evaluate_model(data_train, target_train, data_test, target_test, clf_rf)
253     clf_report_maj = classification_report(target_test, majority_neigh(network, test_idx))
254
255     with open('../results/res_classification.txt', 'w') as f:
256         f.write(clf_report_rf + '\n')
257         f.write(clf_report_maj)
258
259
260 if __name__ == '__main__':
261     main()

```