
You are currently looking at **version 1.2** of this notebook. To download notebooks and datafiles, as well as get help on Jupyter notebooks in the Coursera platform, visit the [Jupyter Notebook FAQ \(https://www.coursera.org/learn/python-social-network-analysis/resources/yPcBs\)](https://www.coursera.org/learn/python-social-network-analysis/resources/yPcBs) course resource.

Assignment 4

```
In [1]: import networkx as nx
import pandas as pd
import numpy as np
import pickle

# import matplotlib.pyplot as plt
# import matplotlib.style as style
# style.use('fivethirtyeight')

from sklearn.linear_model import LogisticRegression

# Hide warnings
import warnings
warnings.filterwarnings('ignore')

# The following lines adjust the granularity of reporting
pd.options.display.max_rows = 10
pd.options.display.float_format = '{:.4f}'.format
```

Part 1 - Random Graph Identification

For the first part of this assignment you will analyze randomly generated graphs and determine which algorithm created them.

```
In [2]: P1_Graphs = pickle.load(open('A4_graphs','rb'))
P1_Graphs

Out[2]: [<networkx.classes.graph.Graph at 0x7fefee786c88>,
<networkx.classes.graph.Graph at 0x7fefee79f080>,
<networkx.classes.graph.Graph at 0x7fefee79f198>,
<networkx.classes.graph.Graph at 0x7fefee79f0b8>,
<networkx.classes.graph.Graph at 0x7fefee79f0f0>]
```

P1_Graphs is a list containing 5 networkx graphs. Each of these graphs were generated by one of three possible algorithms:

- Preferential Attachment ('PA')
- Small World with low probability of rewiring ('SW_L')
- Small World with high probability of rewiring ('SW_H')

Analyze each of the 5 graphs and determine which of the three algorithms generated the graph.

The `graph_identification` function should return a list of length 5 where each element in the list is either 'PA', 'SW_L', or 'SW_H'.

```
In [3]: def degree_distribution(G):
degrees = G.degree()
degree_values = sorted(set(degrees.values()))
histogram = [list(degrees.values()).count(i)/float(nx.number_of_nodes( G)) for i in degree_values]
return histogram
```

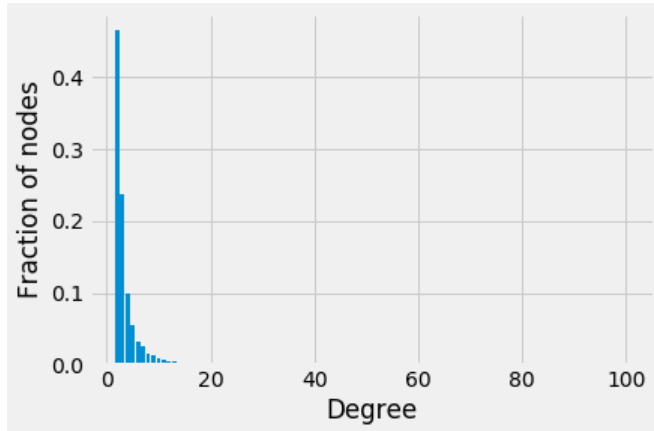
```
In [4]: # G = P1_Graphs[0]

# print(nx.info(G))

# degrees = G.degree()
# degree_values = sorted(set(degrees.values()))
# histogram = [List(degrees.values()).count(i)/float(nx.number_of_nodes(G)) for i in degree_values]

# plt.bar(degree_values, histogram)
# plt.xlabel('Degree')
# plt.ylabel('Fraction of nodes')
# plt.show()
```

Name: barabasi_albert_graph(1000,2)
Type: Graph
Number of nodes: 1000
Number of edges: 1996
Average degree: 3.9920



```
In [5]: def graph_identification():

# Your Code Here
algorithms = []

for G in P1_Graphs:
    clustering = nx.average_clustering(G)
    shortest_path = nx.average_shortest_path_length(G)
    degree_histogram = degree_distribution(G)

    if len(degree_histogram) > 10:
        algorithms.append('PA')

    elif clustering < 0.1:
        algorithms.append('SW_H')

    else:
        algorithms.append('SW_L')

return algorithms

graph_identification()
```

```
Out[5]: ['PA', 'SW_L', 'SW_L', 'PA', 'SW_H']
```

Part 2 - Company Emails

For the second part of this assignment you will be working with a company's email network where each node corresponds to a person at the company, and each edge indicates that at least one email has been sent between two people.

The network also contains the node attributes Department and ManagementSalary.

Department indicates the department in the company which the person belongs to, and ManagementSalary indicates whether that person is receiving a management position salary.

```
In [6]: G = nx.read_gpickle('email_prediction.txt')

print(nx.info(G))
```

```
Name:
Type: Graph
Number of nodes: 1005
Number of edges: 16706
Average degree: 33.2458
```

Part 2A - Salary Prediction

Using network G, identify the people in the network with missing values for the node attribute ManagementSalary and predict whether or not these individuals are receiving a management position salary.

To accomplish this, you will need to create a matrix of node features using networkx, train a sklearn classifier on nodes that have ManagementSalary data, and predict a probability of the node receiving a management salary for nodes where ManagementSalary is missing.

Your predictions will need to be given as the probability that the corresponding employee is receiving a management position salary.

The evaluation metric for this assignment is the Area Under the ROC Curve (AUC).

Your grade will be based on the AUC score computed for your classifier. A model which with an AUC of 0.88 or higher will receive full points, and with an AUC of 0.82 or higher will pass (get 80% of the full points).

Using your trained classifier, return a series of length 252 with the data being the probability of receiving management salary, and the index being the node id.

Example:

```
1      1.0
2      0.0
5      0.8
8      1.0
...
996    0.7
1000   0.5
1001   0.0
Length: 252, dtype: float64
```

```

In [7]: def salary_predictions():

    # Your Code Here
    df = pd.DataFrame(index=G.nodes())

    df['clustering'] = pd.Series(nx.clustering(G))
    df['degree'] = pd.Series(nx.degree(G))
    df['degree_centrality'] = pd.Series(nx.degree_centrality(G))
    df['closeness'] = pd.Series(nx.closeness_centrality(G, normalized=True))
    df['betweenness'] = pd.Series(nx.betweenness_centrality(G, normalized=True))
    df['page_rank'] = pd.Series(nx.pagerank(G))

    df['Department'] = pd.Series(nx.get_node_attributes(G, 'Department'))
    df['ManagementSalary'] = pd.Series(nx.get_node_attributes(G, 'ManagementSalary'))

    train = df[df['ManagementSalary'].notnull()]
    test = df[df['ManagementSalary'].isnull()]

    X_train = train.loc[:, train.columns != 'ManagementSalary']
    y_train = train['ManagementSalary']

    X_test = test.loc[:, test.columns != 'ManagementSalary']
    model = LogisticRegression().fit(X_train, y_train)

    # print('Accuracy of Logistic regression classifier on training set: {:.3f}'\
    #       .format(logisticRegr.score(X_train, y_train)))

    preds = model.predict_proba(X_test)[:, 1]

    probabilities = pd.Series(preds, index = test.index)

    return probabilities

salary_predictions()

```

```

Out[7]: 1      0.1522
        2      0.5895
        5      0.9634
        8      0.1316
       14      0.3551
        ...
      992      0.0159
      994      0.0180
      996      0.0168
     1000      0.0464
     1001      0.0966
Length: 252, dtype: float64

```

Part 2B - New Connections Prediction

For the last part of this assignment, you will predict future connections between employees of the network. The future connections information has been loaded into the variable `future_connections`. The index is a tuple indicating a pair of nodes that currently do not have a connection, and the `Future Connection` column indicates if an edge between those two nodes will exist in the future, where a value of 1.0 indicates a future connection.

```
In [8]: future_connections = pd.read_csv('Future_Connections.csv', index_col=0, converters={0: eval})
future_connections.head(10)
```

Out[8]:

	Future Connection
(6, 840)	0.0000
(4, 197)	0.0000
(620, 979)	0.0000
(519, 872)	0.0000
(382, 423)	0.0000
(97, 226)	1.0000
(349, 905)	0.0000
(429, 860)	0.0000
(309, 989)	0.0000
(468, 880)	0.0000

Using network G and future_connections, identify the edges in future_connections with missing values and predict whether or not these edges will have a future connection.

To accomplish this, you will need to create a matrix of features for the edges found in future_connections using networkx, train a sklearn classifier on those edges in future_connections that have Future Connection data, and predict a probability of the edge being a future connection for those edges in future_connections where Future Connection is missing.

Your predictions will need to be given as the probability of the corresponding edge being a future connection.

The evaluation metric for this assignment is the Area Under the ROC Curve (AUC).

Your grade will be based on the AUC score computed for your classifier. A model which with an AUC of 0.88 or higher will receive full points, and with an AUC of 0.82 or higher will pass (get 80% of the full points).

Using your trained classifier, return a series of length 122112 with the data being the probability of the edge being a future connection, and the index being the edge as represented by a tuple of nodes.

Example:

```
(107, 348)    0.35
(542, 751)    0.40
(20, 426)     0.55
(50, 989)     0.35
...
(939, 940)    0.15
(555, 905)    0.35
(75, 101)     0.65
Length: 122112, dtype: float64
```

```
In [19]: def new_connections_predictions():

    # Your Code Here
    future_connections['preferential_attachment'] = [i[2] for i in nx.preferential_attachment(G, future_connections.index)]

    future_connections['common_neighbors'] = future_connections.index.map(lambda nodes:
                                                                           len(list(nx.common_neighbors(G, nodes[0], nodes[1]
))))))

    train = future_connections[future_connections['Future Connection'].notnull()]
    test = future_connections[future_connections['Future Connection'].isnull()]

    X_train = train.loc[:, train.columns != 'Future Connection']
    y_train = train['Future Connection']

    X_test = test.loc[:, test.columns != 'Future Connection']
    model = LogisticRegression().fit(X_train, y_train)

    preds = model.predict_proba(X_test)[: , 1]

    probabilities = pd.Series(preds, index = test.index)

    return probabilities

new_connections_predictions()
```

```
Out[19]: (107, 348)    0.0377
(542, 751)    0.0141
(20, 426)    0.6450
(50, 989)    0.0144
(942, 986)    0.0147
...
(165, 923)    0.0116
(673, 755)    0.0147
(939, 940)    0.0147
(555, 905)    0.0140
(75, 101)    0.0222
Length: 122112, dtype: float64
```

In []: