

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/vPcBs>) course resource.

---

# Assignment 4

In [1]:

```
import networkx as nx
import pandas as pd
import numpy as np
import pickle
```

## 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 0x7f3ac408f1d0>,
<networkx.classes.graph.Graph at 0x7f3a80a44b70>,
<networkx.classes.graph.Graph at 0x7f3a80a44ba8>,
<networkx.classes.graph.Graph at 0x7f3a80a44be0>,
<networkx.classes.graph.Graph at 0x7f3a80a44c18>]
```

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 graph_identification():

    # Your Code Here
    ans = ['PA', 'SW_L', 'SW_L', 'PA', 'SW_H']

    return ans # Your Answer Here
```

## 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 [4]:

```
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 [5]:

```

from sklearn.linear_model import LogisticRegression

def salary_predictions():

    # Your Code Here
    df = pd.DataFrame(index=G.nodes())
    df['department'] = pd.Series(nx.get_node_attributes(G, 'Department'))
    df['clustering'] = pd.Series(nx.clustering(G))
    df['degree'] = pd.Series(G.degree())
    df['salary'] = pd.Series(nx.get_node_attributes(G, 'ManagementSalary'))

    df1 = df.dropna(axis=0, how='any')
    df2 = df[df['salary'].isnull()]

    X_train = df1[['department', 'clustering', 'degree']]
    y_train = df1[['salary']]
    X_test = df2[['department', 'clustering', 'degree']]

    lr = LogisticRegression().fit(X_train, y_train)
    lr_predicted = lr.predict_proba(X_test)

    ans = pd.Series(data = lr_predicted[:,1], index = df2.index)

    return ans # Your Answer Here

```

## 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 [6]:

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

Out[6]:

**Future Connection**

(6, 840)	0.0
(4, 197)	0.0
(620, 979)	0.0
(519, 872)	0.0
(382, 423)	0.0
(97, 226)	1.0
(349, 905)	0.0
(429, 860)	0.0
(309, 989)	0.0
(468, 880)	0.0

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 [7]:

```
def new_connections_predictions():

    # Your Code Here
    future_connections['preferential attachment'] = [i[2] for i in nx.preferential_attachment_index(G)]
    future_connections['Common emp'] = future_connections.index.map(lambda emp: len(list(nx.common_neighbors(G, emp, emp))))
    future_connections['jaccard'] = [i[2] for i in nx.jaccard_coefficient(G, future_connections.index, future_connections.index)]
    future_connections['rs alloc'] = [i[2] for i in nx.resource_allocation_index(G, future_connections.index, future_connections.index)]
    future_connections['adar'] = [i[2] for i in nx.adamic_adar_index(G, future_connections.index, future_connections.index)]

    df1 = future_connections.dropna(axis=0, how='any')
    df2 = future_connections[future_connections['Future Connection'].isnull()]

    X_train = df1[['preferential attachment', 'Common emp', 'jaccard', 'rs alloc', 'adar']]
    y_train = df1[['Future Connection']]
    X_test = df2[['preferential attachment', 'Common emp', 'jaccard', 'rs alloc', 'adar']]

    lr = LogisticRegression().fit(X_train, y_train)
    lr_predicted = lr.predict_proba(X_test)

    ans = pd.Series(data = lr_predicted[:,1], index = df2.index)

    return ans # Your Answer Here
```