

Clustering Assignment

There will be some functions that start with the word "grader" ex: grader_actors(), grader_movies(), grader_cost1() etc, you should not change those function definition.

Every Grader function has to return True.
Please check [clustering assignment helper functions](#) notebook before attempting this assignment.

- Read graph from the given `movie_actor_network.csv` (note that the graph is bipartite graph)
- Using stellergaph and gensim packages, get the dense representation (128dimensional vector) of every node in the graph. [Refer [Clustering_Assignment_Reference.ipynb](#)]
- Split the dense representation into actor nodes, movies nodes.(Write you code in `def data_split()`)

Task 1 : Apply clustering algorithm to group similar actors

- For this task consider only the actor nodes
- Apply any clustering algorithm of your choice
Refer : <https://scikit-learn.org/stable/modules/clustering.html>
- Choose the number of clusters for which you have maximum score of $C_{cost1} + C_{cost2}$
- $Cost1 = \frac{1}{N} \sum_{i=1}^N \sum_{j \in cluster i} (number\ of\ nodes\ in\ the\ largest\ connected\ component\ in\ the\ graph\ with\ the\ actor\ node\ and\ its\ movie\ neighbours\ in\ cluster\ i)$ where N= number of clusters
(Write your code in `def cost1()`)
(total number of nodes in that cluster i)
- $Cost2 = \frac{1}{N} \sum_{i=1}^N \sum_{j \in cluster i} (sum\ of\ degrees\ of\ actor\ nodes\ in\ the\ graph\ with\ the\ actor\ node\ and\ its\ movie\ neighbours\ in\ cluster\ i)$ where N= number of clusters
(Write your code in `def cost2()`)
(number of unique movie nodes in the graph with the actor node and its movie neighbours in cluster i)
- Fit the clustering algorithm with the optimal number_of_clusters and get the cluster number for each node
- Convert the d-dimensional dense vectors of nodes into 2-dimensional using dimensionality reduction techniques (preferably TSNE)
- Plot the 2d scatter plot, with the node vectors after step 6 and give colors to nodes such that same cluster nodes will have same color

Task 2 : Apply clustering algorithm to group similar movies

- For this task consider only the movie nodes
- Apply any clustering algorithm of your choice 3 Choose the number of clusters for which you have maximum score of $C_{cost1} + C_{cost2}$
- $Cost1 = \frac{1}{N} \sum_{i=1}^N \sum_{j \in cluster i} (number\ of\ nodes\ in\ the\ largest\ connected\ component\ in\ the\ graph\ with\ the\ movie\ node\ and\ its\ actor\ neighbours\ in\ cluster\ i)$ where N= number of clusters
(Write your code in `def cost1()`)
(total number of nodes in that cluster i)
- $Cost2 = \frac{1}{N} \sum_{i=1}^N \sum_{j \in cluster i} (sum\ of\ degrees\ of\ movie\ nodes\ in\ the\ graph\ with\ the\ movie\ node\ and\ its\ actor\ neighbours\ in\ cluster\ i)$ where N= number of clusters
(Write your code in `def cost2()`)
(number of unique actor nodes in the graph with the movie node and its actor neighbours in cluster i)

Algorithm for actor nodes

```
for number_of_clusters in [3, 5, 10, 30, 50, 100, 200, 500]:
    algo = clustering_algorithm(clusters=number_of_clusters)
    # you will be passing a matrix of size N*d where N number of actor nodes and d is dimension from gensim
    algo.fit(the dense vectors of actor nodes)
    # You can get the labels for corresponding actor nodes (algo.labels_)
    # Create a graph for every cluster (ie., if n_clusters=3, create 3 graphs)
    # You can use ego_graph to create subgraph from the actual graph
    compute cost1,cost2
    # (if n_clusters=3, cost1=cost1(graph1)+cost1(graph2)+cost1(graph3) # here we are doing summation
    cost2=cost2(graph1)+cost2(graph2)+cost2(graph3)
    compute the metric Cost = Cost1+Cost2
    return number_of_clusters which have maximum Cost
```

```
In [308]: pip install networkx==2.5.1
Requirement already satisfied: networkx==2.5.1 in c:\users\subhashini rajesh\anaconda3\lib\site-packages (2.5.1)
Requirement already satisfied: decorator<4.3,=>4.3 in c:\users\subhashini rajesh\anaconda3\lib\site-packages (from networkx==2.5.1) (4.4.2)
```

```
In [1]: import networkx as nx
from networkx.algorithms import bipartite
import matplotlib.pyplot as plt
from sklearn.cluster import KMeans
import numpy as np
import warnings
warnings.filterwarnings("ignore")
import pandas as pd
# you need to have tensorflow
from stellargraph.data import UniformRandomMetaPathWalk
from stellargraph import StellarGraph
from collections import defaultdict

In [2]: nx.__version__
Out[2]: '2.5.1'

In [5]: import matplotlib
matplotlib.__version__
Out[5]: '3.5.1'

In [3]: datapd.read_csv('C:\Users\Subhashini Rajesh\Graph model AAIC\movie_actor_network.csv', index_col=False, names=['movie', 'actor'])

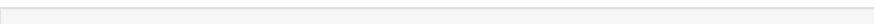
In [4]: edges = [(tuple(x) for x in data.values.tolist())]

In [4]: B = nx.Graph()
B.add_nodes_from(data['movie'].unique(), bipartite=0, label='movie')
B.add_nodes_from(data['actor'].unique(), bipartite=1, label='actor')
B.add_edges_from(edges, label='acted')
```

```
In [5]: A = list(B.subgraph(c) for c in nx.connected_components(B))[0]

In [6]: print("number of nodes", A.number_of_nodes())
print("number of edges", A.number_of_edges())
number of nodes 4783
number of edges 9650

In [7]: l, r = nx.bipartite.sets(A)
pos = {}
pos.update((node, (1, index)) for index, node in enumerate(l))
pos.update((node, (2, index)) for index, node in enumerate(r))
nx.draw(A, pospos, with_labels=True)
plt.show()
```



```
In [9]: # Create the random walker
rw = UniformRandomMetaPathWalk(StellarGraph(A))

# specify the metapath schemas as a list of lists of node types.
metapaths = [
    ["movie", "actor", "movie"],
    ["actor", "movie", "actor"]
]

walks = rw.run(nodes=list(A.nodes()), # root nodes
               length=100, # maximum length of a random walk
               nts, # number of random walks per root node
               metapaths=metapaths
            )
print("Number of random walks: {}".format(len(walks)))
Number of random walks: 4703
```

```
In [10]: from gensim.models import Word2Vec
model = Word2Vec(walks, vector_size=128, window=5)
```

```
In [11]: model.wv.vectors.shape # 128-dimensional vector for each node in the graph
Out[11]: (4703, 128)
```

```
In [12]: # Retrieve node embeddings and corresponding subjects
node_ids = model.index_to_key # List of node IDs
node_embeddings = model.wv.vectors # numpy.ndarray of size number of nodes times embeddings dimensionality
node_targets = [A.nodes[node_id]['label'] for node_id in node_ids]
node_embeddings.shape
```

```
Out[12]: (4703, 128)
```

```
print(node_ids[:15], end='')

['a973', 'a967', 'a964', 'a1731', 'a960', 'a970', 'a1028', 'a1057', 'a965', 'a1003', 'm1094', 'a966', 'm67', 'a988', 'm1111']
```

```
print(node_targets[:15],end='')

['actor', 'actor', 'actor', 'actor', 'actor', 'actor', 'actor', 'actor', 'actor', 'movie', 'actor', 'movie', 'actor', 'movie']
```

```
In [13]: len(node_embeddings)
Out[13]: 4703
```

```
In [14]: def data_split(node_ids,node_targets,node_embeddings):
'''In this function, we will split the data into actor_embeddings , movie_embeddings '''
actor_nodes,movie_nodes=[],[]
actor_embeddings,movie_embeddings=[],[]
actor_ids,movie_ids = [],[]
# split the node_embeddings into actor_embeddings,movie_embeddings based on node_ids
# By using node_embedding and node_targets, we can extract actor_embedding and movie_embedding
# By using node_ids and node_targets, we can extract actor_nodes and movie nodes
d1 = {node_ids[i] : node_embeddings[i] for i in range(len(node_ids))}
d2 = {node_ids[i] : node_targets[i] for i in range(len(node_ids))}
for k in d1.keys():
    if j[0] == 'a':
        actor_embeddings.append(d1[j])
        actor_ids.append(j)
    if j[0] == 'm':
        movie_embeddings.append(d1[j])
        movie_ids.append(j)
    for i in d2.values():
        if i == 'actor':
            actor_nodes.append(i)
        else:
            movie_nodes.append(i)
return actor_nodes,movie_nodes,actor_embeddings,movie_embeddings,actor_ids,movie_ids
```

```
In [15]: actor_nodes,movie_nodes,actor_embeddings,movie_embeddings,actor_ids,movie_ids = data_split(node_ids,node_targets,node_embeddings)
actor_embeddings = np.array(actor_embeddings)
movie_embeddings = np.array(movie_embeddings)
print(movie_embeddings.shape)
print(actor_embeddings.shape)
```

```
Out[15]: (1292, 128)
(3411, 128)
```

```
In [16]: print(len(actor_ids))
3411
```

```
Grader function -1
def grader_actors(data):
assert(len(data)==3411)
return True
grader_actors(actor_nodes)
```

```
Out[17]: True
Grader function -2
```

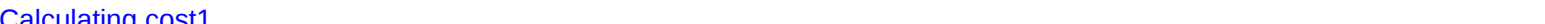
```
In [18]: def grader_movies(data):
assert(len(data)==1292)
return True
grader_movies(movie_nodes)
```

```
Out[18]: True
Calculating cost1
```

$Cost1 = \frac{1}{N} \sum_{i=1}^N \sum_{j \in cluster i} (number\ of\ nodes\ in\ the\ largest\ connected\ component\ in\ the\ graph\ with\ the\ actor\ node\ and\ its\ movie\ neighbours\ in\ cluster\ i)$ where N= number of clusters
(total number of nodes in that cluster i)

```
In [19]: def cost1(graph,number_of_clusters):
'''In this function, we will calculate cost1'''
total_number_of_nodes = graph.number_of_nodes()
graphs = [graph.subgraph(c) for i in range(nx.connected_components(graph))
max_nodes = max([len(i.nodes()) for i in graphs])
cost1 = ((number_of_clusters) * (max_nodes*total_number_of_nodes))
return cost1
```

```
In [20]: import networkx as nx
from networkx.algorithms import bipartite
graded_graph= nx.Graph()
graded_graph.add_nodes_from(['a1','a5','a10','a11'], bipartite=0) # Add the node attribute "bipartite"
graded_graph.add_nodes_from(['m1','m2','m4','m6','m5','m8'], bipartite=1)
graded_graph.add_edges_from([(('a1','m1'),('a1','m2'),('a1','m4'),('a1','m6'),('a1','m5'),('a1','m8')],
l=['a1','a5','a10','a11'];r=['m1','m4','m6','m5','m8'])
pos = {}
pos.update((node, (1, index)) for index, node in enumerate(l))
pos.update((node, (2, index)) for index, node in enumerate(r))
nx.draw_networkx(graded_graph, pospos, with_labels=True,node_color='lightblue',alpha=0.8,style='dotted',node_size=500)
```



```
In [21]: # print(set(data['movie']))
Grader function -3
```

```
In [22]: graded_cost1(cost1(graded_graph,3)
def grader_cost1(data):
assert(data==(1/3)*(4/10)) # 1/3 is number of clusters
return True
grader_cost1(graded_cost1)
```

```
Out[22]: True
Calculating cost2
```

$Cost2 = \frac{1}{N} \sum_{i=1}^N \sum_{j \in cluster i} (sum\ of\ degrees\ of\ actor\ nodes\ in\ the\ graph\ with\ the\ actor\ node\ and\ its\ movie\ neighbours\ in\ cluster\ i)$ where N= number of clusters
(total number of nodes in that cluster i)

```
In [23]: def cost2(graph,number_of_clusters):
'''In this function, we will calculate cost1'''
degree = graph.degree(data['actor'])
degree = dict(degree)
deg_sum = sum(degree.values())
movies = []
actors = []
for i in graph.nodes():
    if 'a' in i:
        movies.append(i)
    if 'a' in i:
        actors.append(i)
cost2 = ((1/number_of_clusters)*(deg_sum/len(movies)))
return cost2
```

```
In [24]: def cost2_mov(graph,number_of_clusters):
'''In this function, we will calculate cost1'''
degree = graph.degree(data['actor'])
degree = dict(degree)
deg_sum = sum(degree.values())
movies = []
actors = []
for i in graph.nodes():
    if 'a' in i:
        movies.append(i)
    if 'a' in i:
        actors.append(i)
cost2 = ((1/number_of_clusters)*(deg_sum/len(actors)))
return cost2
```

```
Grader function -4
graded_cost2=cost2(graded_graph,3)
def grader_cost2(data):
assert(data==(1/3)*(6/6)) # 1/3 is number of clusters
return True
grader_cost2(graded_cost2)
```

```
Out[25]: True
Grouping similar actors
```

```
In [26]: number_of_clusters = [3, 5, 10, 30, 50, 100, 200, 500]
cost_dict = {}
for no_of_cluster in number_of_clusters:
    cost_1=0
    cost_2=0
    model = KMeans(no_of_cluster)
    model.fit(actor_embeddings)
    labels = model.labels_
    d1 = {actor_ids[i]:labels[i] for i in range(len(actor_ids))}
    cls = defaultdict(list) #default dictionary to frame clusters and its actor nodes
    for k,v in d1.items(): # example ('a1':3, 'a2':2, 'a3':2, 'a4':3) num_of_cluster = 3
        cls[v].append(k) # (1: ['a1'], 2: ['a2'], 'a3'], 3: ['a4'])
    for k, v in cls.items():
        G = nx.Graph()
        total_cost = 0
        for i in range(len(v)): # iterating actor nodes
            eg = nx.ego_graph(B, v[i])
            G.add_nodes_from(eg.nodes())
            G.add_edges_from(eg.edges())
            # print(cost2(G, cluster))
            # print(cost1(G, cluster))
            cost_1 += cost1(G, no_of_cluster)
            cost_2 += cost2(G, no_of_cluster)
        # print(cost_1, cost_2)
        total_cost = cost_1 + cost_2
        # print(total_cost)
    cost_dict[no_of_cluster] = total_cost
```

```
In [27]: cost_dict
Out[27]: {3: 3.718090842947067,
5: 3.017589919361510,
10: 2.284958067723567,
30: 1.7538954249587,
50: 1.520525130320901,
100: 1.3528799323694219,
200: 1.027623234485072,
500: 1.062216327604064}
```

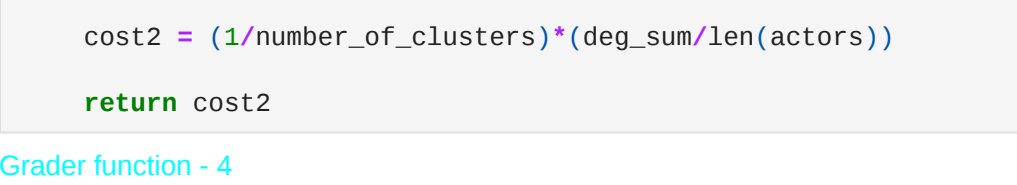
Displaying similar actor clusters

```
In [28]: #optimized k for kmeans is 3
#max of cost1 + cost2 from number_of_clusters, optimized k is 3
opt_cluster = max(zip(cost_dict.values(), cost_dict.keys()))[1]
model = KMeans(opt_cluster)
model.fit(actor_embeddings)
labels_actor = model.labels_
```

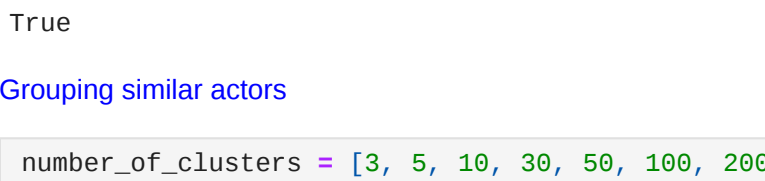
```
In [29]: from sklearn.manifold import TSNE
transform = TSNE #PCA
trans = transform(n_components=2,perplexity=50)
node_embeddings_2d = trans.fit_transform(actor_embeddings)
```

```
Out[34]: node_embeddings_2d[:1]
array([[ -26.962755,   55.62998 ],
       [ -25.326025,   51.051733 ],
       [ -20.060207,   55.295824 ]],
      dtype=float32)
```

```
In [35]: import numpy as np
# draw the points
label_map = { i: i for i, i in enumerate(np.unique(labels_actor))}
node_colours = [label_map[target] for target in labels_actor]
plt.figure(figsize=(10,8))
plt.axes().set(aspect='equal')
plt.scatter(node_embeddings_2d[:,0],
            node_embeddings_2d[:,1],
            c=node_colours, alpha=0.3)
plt.title('() visualization of node embeddings'.format(transform.__name__))
plt.show()
```



```
In [36]: import matplotlib.pyplot as plt
import seaborn as sns
plt.scatter(actor_ids, labels_actor, label = [1,2,3])
plt.xlabel('Actor Nodes')
plt.ylabel('Clusters')
plt.show()
```



Grouping similar movies

```
In [37]: number_of_clusters = [3, 5, 10, 30, 50, 100, 200, 500]
cost_dict_mov = {}
for no_of_cluster in number_of_clusters:
    cost_1=0
    cost_2=0
    model = KMeans(no_of_cluster)
    model.fit(movie_embeddings)
    labels = model.labels_
    d1 = {movie_ids[i]:labels[i] for i in range(len(movie_ids))}
    cls = defaultdict(list) #default dictionary to frame clusters and its actor nodes
    for k,v in d1.items(): # example ('a1':3, 'a2':2, 'a3':2, 'a4':3) num_of_cluster = 3
        cls[v].append(k) # (1: ['a1'], 2: ['a2'], 'a3'], 3: ['a4'])
    for k, v in cls.items():
        G = nx.Graph()
        total_cost = 0
        for i in range(len(v)): # iterating actor nodes
            eg = nx.ego_graph(B, v[i])
            G.add_nodes_from(eg.nodes())
            G.add_edges_from(eg.edges())
            # print(cost2(G, cluster))
            # print(cost1(G, cluster))
            cost_1 += cost1(G, no_of_cluster)
            cost_2 += cost2_mov(G, no_of_cluster)
        # print(cost_1, cost_2)
        total_cost = cost_1 + cost_2
        # print(total_cost)
    cost_dict_mov[no_of_cluster] = total_cost
```

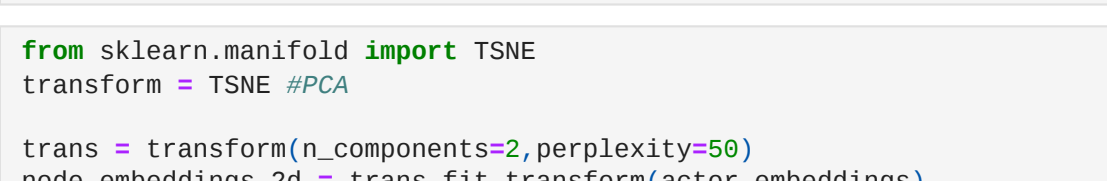
```
In [56]: cost_dict_mov
Out[56]: {3: 3.2046369623516,
5: 2.960295716386205,
10: 2.437266024508065,
30: 2.0417087251771024,
50: 1.82326355981011,
100: 1.625430938962629,
200: 1.377878448780635,
500: 1.204135232448585}
```

Displaying similar movie clusters

```
In [38]: #optimized for kmeans is 3
#max of cost1 + cost2 from number_of_clusters, optimized k is 3
opt_cluster = max(zip(cost_dict_mov.values(), cost_dict.keys()))[1]
model = KMeans(opt_cluster)
model.fit(movie_embeddings)
labels_movie = model.labels_
```

```
In [39]: from sklearn.manifold import TSNE
transform = TSNE #PCA
trans = transform(n_components=2)
node_embeddings_2d = trans.fit_transform(movie_embeddings)
```

```
In [40]: import numpy as np
# draw the points
label_map = { i: i for i, i in enumerate(np.unique(labels_movie))}
node_colours = [label_map[target] for target in labels_movie]
plt.figure(figsize=(10,8))
plt.axes().set(aspect='equal')
plt.scatter(node_embeddings_2d[:,0],
            node_embeddings_2d[:,1],
            c=node_colours, alpha=0.3)
plt.title('() visualization of node embeddings'.format(transform.__name__))
plt.show()
```



```
In [42]: import matplotlib.pyplot as plt
import seaborn as sns
plt.scatter(movie_ids, labels_movie, label = [1,2,3])
plt.xlabel('Movie Nodes')
plt.ylabel('Clusters')
plt.show()
```

