

# 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.

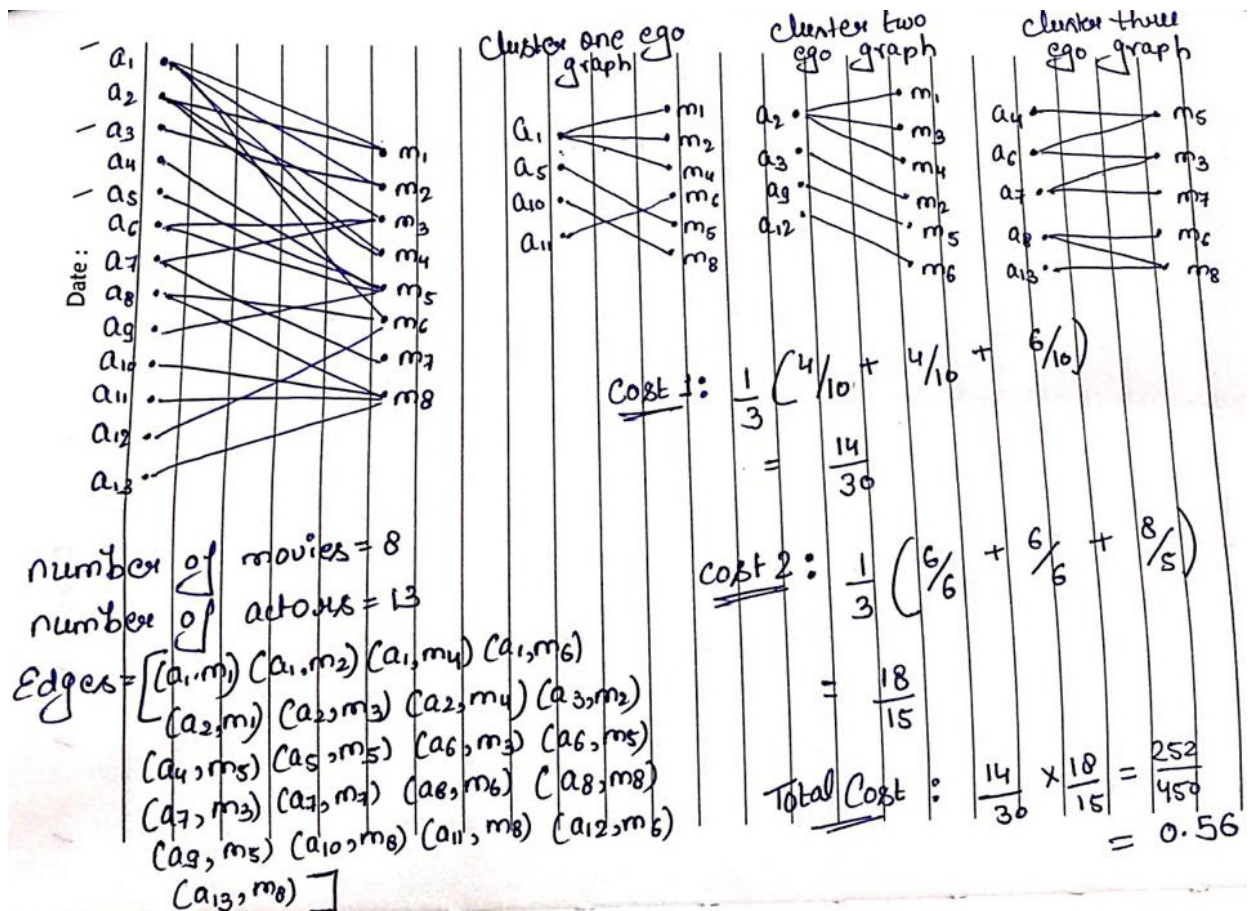
Please check [clustering assignment helper functions](https://drive.google.com/file/d/1V29KhKo3YnckMX32treEgdtH5r90DljU/view?usp=sharing) (<https://drive.google.com/file/d/1V29KhKo3YnckMX32treEgdtH5r90DljU/view?usp=sharing>) notebook before attempting this assignment.

- 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

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





## Task 2 : Apply clustering algorithm to group similar movies

1. For this task consider only the movie nodes
2. Apply any clustering algorithm of your choice 3. Choose the number of clusters for which you have maximum score of  $Cost1 * Cost2$

Cost1 =

$$\frac{1}{N} \sum_{\text{each cluster } i} \frac{(\text{number of nodes in the largest connected component in the graph with the movie nodes and its actor neighbours})}{(\text{total number of nodes in that cluster } i)}$$

where N= number of clusters

(Write your code in `def cost1()`)

3. Cost2 =

$$\frac{1}{N} \sum_{\text{each cluster } i} \frac{(\text{sum of degrees of movie nodes in the graph with the movie nodes and its actor neighbours in cluster } i)}{(\text{number of unique actor nodes in the graph with the movie nodes and its actor neighbours in cluster } i)}$$

where N= number of clusters

(Write your code in `def cost2()`)

### 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_cluster=3, cost1=cost1(graph1)+cost1(graph2)+cost1(graph
3) # 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 [1]: 1 import networkx as nx
2 from networkx.algorithms import bipartite
3 import matplotlib.pyplot as plt
4 from sklearn.cluster import KMeans
5 import numpy as np
6 import warnings
7 warnings.filterwarnings("ignore")
8 import pandas as pd
9 # you need to have tensorflow
10 from stellargraph.data import UniformRandomMetaPathWalk
11 from stellargraph import StellarGraph
12 from tqdm import tqdm

```

```

In [2]: 1 data=pd.read_csv('movie_actor_network.csv', index_col=False, names=['movie', '
2

```

```

In [3]: 1 edges = [tuple(x) for x in data.values.tolist()]

```

```

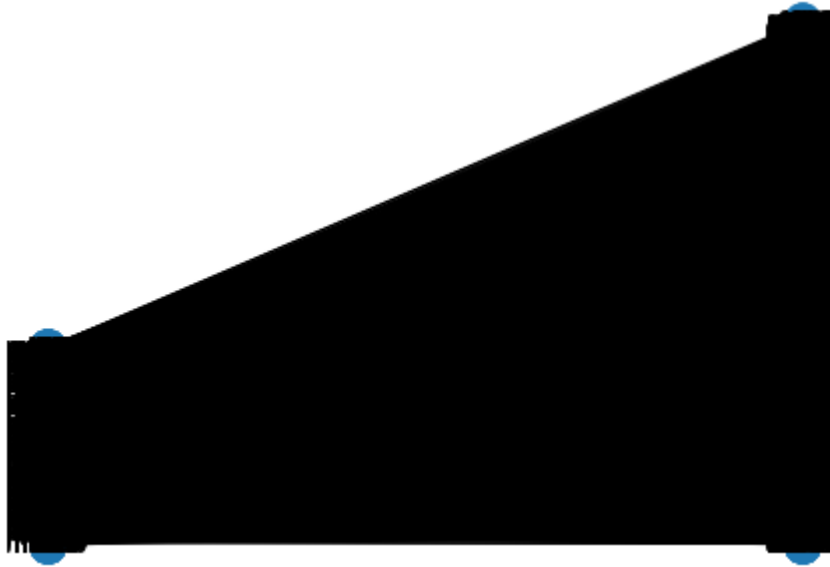
In [4]: 1 edges=[(x,y) for (y,x) in edges]

```

```
In [5]: 1 B = nx.Graph()
2 B.add_nodes_from(data['movie'].unique(), bipartite=0, label='movie')
3 B.add_nodes_from(data['actor'].unique(), bipartite=1, label='actor')
4 B.add_edges_from(edges, label='acted')
5 A = list(nx.connected_component_subgraphs(B))[0]
6 print("number of nodes", A.number_of_nodes())
7 print("number of edges", A.number_of_edges())
8 l, r = nx.bipartite.sets(A)
9 pos = {}
10
11 pos.update((node, (1, index)) for index, node in enumerate(l))
12 pos.update((node, (2, index)) for index, node in enumerate(r))
13
14 nx.draw(A, pos=pos, with_labels=True)
15 plt.show()
```

number of nodes 4703

number of edges 9650



```
In [6]: 1 movies = []
2 actors = []
3 for i in B.nodes():
4     if 'm' in i:
5         movies.append(i)
6     if 'a' in i:
7         actors.append(i)
8 print('number of movies ', len(movies))
9 print('number of actors ', len(actors))
```

number of movies 1292  
number of actors 3411

```
In [7]: 1 # Create the random walker
2 rw = UniformRandomMetaPathWalk(StellarGraph(B))
3
4 # specify the metapath schemas as a list of lists of node types.
5 metapaths = [
6     ["movie", "actor", "movie"],
7     ["actor", "movie", "actor"]
8 ]
9
10
11 walks = rw.run(nodes=list(B.nodes()), # root nodes
12               length=100, # maximum length of a random walk
13               n=1, # number of random walks per root node
14               metapaths=metapaths
15               )
16
17 print("Number of random walks: {}".format(len(walks)))
```

Number of random walks: 4703

```
In [8]: 1 from gensim.models import Word2Vec
2 model = Word2Vec(walks, vector_size=128, window=5)
```

```
In [9]: 1 model.wv.vectors.shape # 128-dimensional vector for each node in the graph
```

Out[9]: (4703, 128)

```
In [10]: 1 # Retrieve node embeddings and corresponding subjects
2 node_ids = model.wv.index_to_key # list of node IDs
3 node_embeddings = model.wv.vectors # numpy.ndarray of size number of nodes t
```

```
In [11]: 1 actor_nodes=[i for i in node_ids if "a" in i]
2 movie_nodes=[j for j in node_ids if "m" in j]
3 len(actor_nodes),len(movie_nodes)
```

Out[11]: (3411, 1292)

```

In [12]: 1 def data_split(node_ids,node_embeddings):
2         '''In this function, we will split the node embeddings into actor_embeddings
3         '''First we have taken the node_ids from wv embeddings then tried to get
4         Then with the help of these indices we can find the word embeddings of the
5         nodes_actor=[]
6         nodes_movie=[]
7         for actor_movie in node_ids:
8             if "a" in actor_movie:
9                 nodes_actor.append(actor_movie)
10            else:
11                nodes_movie.append(actor_movie)
12        movie_indices=[]
13        actor_indices=[]
14        for actormovie_indices in node_ids:
15            if "m" in actormovie_indices:
16                movie_indices.append(node_ids.index(actormovie_indices))
17            else:
18                actor_indices.append(node_ids.index(actormovie_indices))
19        actor_embeddings,movie_embeddings=node_embeddings[actor_indices,:],node_embeddings[movie_indices,:]
20        return nodes_actor,nodes_movie,actor_embeddings,movie_embeddings

```

#### Grader function - 1

```

In [13]: 1 def grader_actors(data):
2         assert(len(data)==3411)
3         return True
4         grader_actors(actor_nodes)

```

Out[13]: True

#### Grader function - 2

```

In [14]: 1 def grader_movies(data):
2         assert(len(data)==1292)
3         return True
4         grader_movies(movie_nodes)

```

Out[14]: True

#### Calculating cost1

Cost1 =

$$\frac{1}{N} \sum_{\text{each cluster } i} \frac{(\text{number of nodes in the largest connected component in the graph with the actor nodes and its movie neighbors})}{(\text{total number of nodes in that cluster } i)}$$

where N= number of clusters



```
In [15]: 1 def cost1(graph,number_of_clusters):
2         '''In this function, we will calculate cost1'''
3         max_nodes,total_nodes=len(max(nx.connected_components(graph),key=len)),gr
4         cost1= max_nodes/(total_nodes)
5         return cost1/number_of_clusters
```

### Calculating cost2

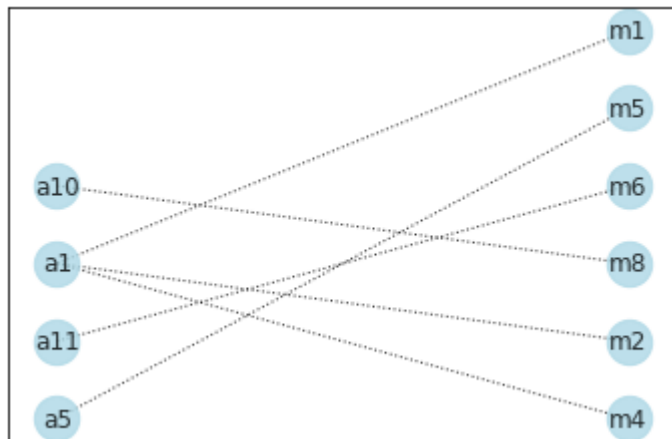
Cost2 =

$$\frac{1}{N} \sum_{\text{each cluster } i} \frac{(\text{sum of degrees of actor nodes in the graph with the actor nodes and its movie neighbours in cluster } i)}{(\text{number of unique movie nodes in the graph with the actor nodes and its movie neighbours in cluster } i)}$$

where N= number of clusters

```
In [16]: 1 def cost2(graph,number_of_clusters):
2         '''In this function, we will calculate cost2'''
3         num=0
4         den=0
5         for i in graph.nodes:
6             if "a" in i:
7                 num+=graph.degree(i)
8             else:
9                 den+=1
10        cost2= num/(den)
11        return cost2/number_of_clusters
```

```
In [17]: 1 import networkx as nx
2         from networkx.algorithms import bipartite
3         graded_graph= nx.Graph()
4         graded_graph.add_nodes_from(['a1','a5','a10','a11'], bipartite=0) # Add the n
5         graded_graph.add_nodes_from(['m1','m2','m4','m6','m5','m8'], bipartite=1)
6         graded_graph.add_edges_from([('a1','m1'),('a1','m2'),('a1','m4'),('a11','m6')
7         l={'a1','a5','a10','a11'};r={'m1','m2','m4','m6','m5','m8'}
8         pos = {} # it is basically required by the nx.draw to set the position of the
9         pos.update((node, (1, index)) for index, node in enumerate(l))
10        pos.update((node, (2, index)) for index, node in enumerate(r))
11        nx.draw_networkx(graded_graph, pos=pos, with_labels=True,node_color='lightblue')
```



## Grader function - 3

```
In [18]: 1 graded_cost1=cost1(graded_graph,3)
2 def grader_cost1(data):
3     assert(data==((1/3)*(4/10))) # 1/3 is number of clusters
4     return True
5 grader_cost1(graded_cost1)
```

Out[18]: True

## Grader function - 4

```
In [19]: 1 graded_cost2=cost2(graded_graph,3)
2 def grader_cost2(data):
3     assert(data==((1/3)*(6/6))) # 1/3 is number of clusters
4     return True
5 grader_cost2(graded_cost2)
```

Out[19]: True

```
In [20]: 1 def Kmeans_clustering(nodes,list_num_cluster,embeddings):
2     all_cost=[]
3     #Running loop for different number of cluster values
4     for k in (list_num_cluster):
5         kmeans = KMeans(n_clusters=k).fit(embeddings)
6         currcost1=0
7         currcost2=0
8         #Finding cost function for each clusters formed from the above fitting
9         for i in range(0,k):
10            current_nodes=np.array(nodes).reshape(len(nodes),)[kmeans.labels_
11            current_graph=nx.Graph()
12            #Creating the graph Cluster
13            for j in current_nodes:
14                current_graph.add_nodes_from(nx.ego_graph(B,j).nodes)
15                current_graph.add_edges_from(nx.ego_graph(B,j).edges())
16                currcost1=currcost1+cost1(current_graph,k)
17                currcost2=currcost2+cost2(current_graph,k)
18            all_cost.append(currcost1*currcost2)
19            kmax=list_num_cluster[all_cost.index(max(all_cost))]
20            return kmax,all_cost
```

```
In [21]: 1 actor_nodes,movie_nodes,actor_embeddings,movie_embeddings=data_split(node_ids
```



```
In [22]: 1 list_num_cluster=[3, 5, 10, 30, 50, 100, 200, 500]
2 kmax_actor,score_actor=Kmeans_clustering(actor_nodes,list_num_cluster,actor_e
3 for i,j in zip(list_num_cluster,score_actor):
4     print("The cost of ",i," number of cluster is ",j)
5     print("The value of k is ",kmax_actor)
```

```
The cost of 3 number of cluster is 3.725575244657303
The cost of 5 number of cluster is 2.888584960475906
The cost of 10 number of cluster is 2.2689389916515634
The cost of 30 number of cluster is 1.7574169111257325
The cost of 50 number of cluster is 1.5670502480239719
The cost of 100 number of cluster is 1.5592628874662784
The cost of 200 number of cluster is 1.8541525553163092
The cost of 500 number of cluster is 1.8470806873477434
The value of k is 3
```

```
In [23]: 1 list_num_cluster=[3, 5, 10, 30, 50, 100, 200, 500]
2 kmax_movie,score_movie=Kmeans_clustering(movie_nodes,list_num_cluster,movie_e
3 for i,j in zip(list_num_cluster,score_movie):
4     print("The cost of ",i," number of cluster is ",j)
5     print("The value of k is ",kmax_movie)
```

```
The cost of 3 number of cluster is 8.388011950992288
The cost of 5 number of cluster is 8.16310541178454
The cost of 10 number of cluster is 8.964492415785722
The cost of 30 number of cluster is 12.72315556651753
The cost of 50 number of cluster is 14.725889907703424
The cost of 100 number of cluster is 13.99838405964806
The cost of 200 number of cluster is 12.88256578652694
The cost of 500 number of cluster is 10.32299746739229
The value of k is 50
```

```
In [24]: 1 kmeans = KMeans(n_clusters=kmax_actor, random_state=0).fit(actor_embeddings)
```

```
In [25]: 1 from sklearn.manifold import TSNE
2 tsne = TSNE(n_components=2,n_jobs=-1)
3 actor_tsne_2d = tsne.fit_transform(actor_embeddings)
```

```
In [26]: 1 #here we are
2 label_map = { l: i for i, l in enumerate(np.unique(list(kmeans.labels_)))}
3 node_colours = [ label_map[target] for target in list(kmeans.labels_)]
4 plt.figure(figsize=(12,12))
5 plt.axes().set(aspect="equal")
6 plt.scatter(actor_tsne_2d[:,0],
7             actor_tsne_2d[:,1],
8             c=node_colours, alpha=1)
9 plt.title('Tsne visualization of node embeddings')
10 plt.show()
```



```
In [27]: 1 kmeans = KMeans(n_clusters=kmax_movie, random_state=0).fit(movie_embeddings)
```

```
In [28]: 1 from sklearn.manifold import TSNE
2 tsne = TSNE(n_components=2,n_jobs=-1)
3 movie_tsne_2d = tsne.fit_transform(movie_embeddings)
```

```
In [29]: 1 label_map = { l: i for i, l in enumerate(np.unique(list(kmeans.labels_)))}
2 node_colours = [ label_map[target] for target in list(kmeans.labels_)]
3 plt.figure(figsize=(12,12))
4 plt.axes().set(aspect="equal")
5 plt.scatter(movie_tsne_2d[:,0],
6             movie_tsne_2d[:,1],
7             c=node_colours, alpha=1)
8 plt.title('tsne visualization of node embeddings')
9 plt.show()
```