# 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.**
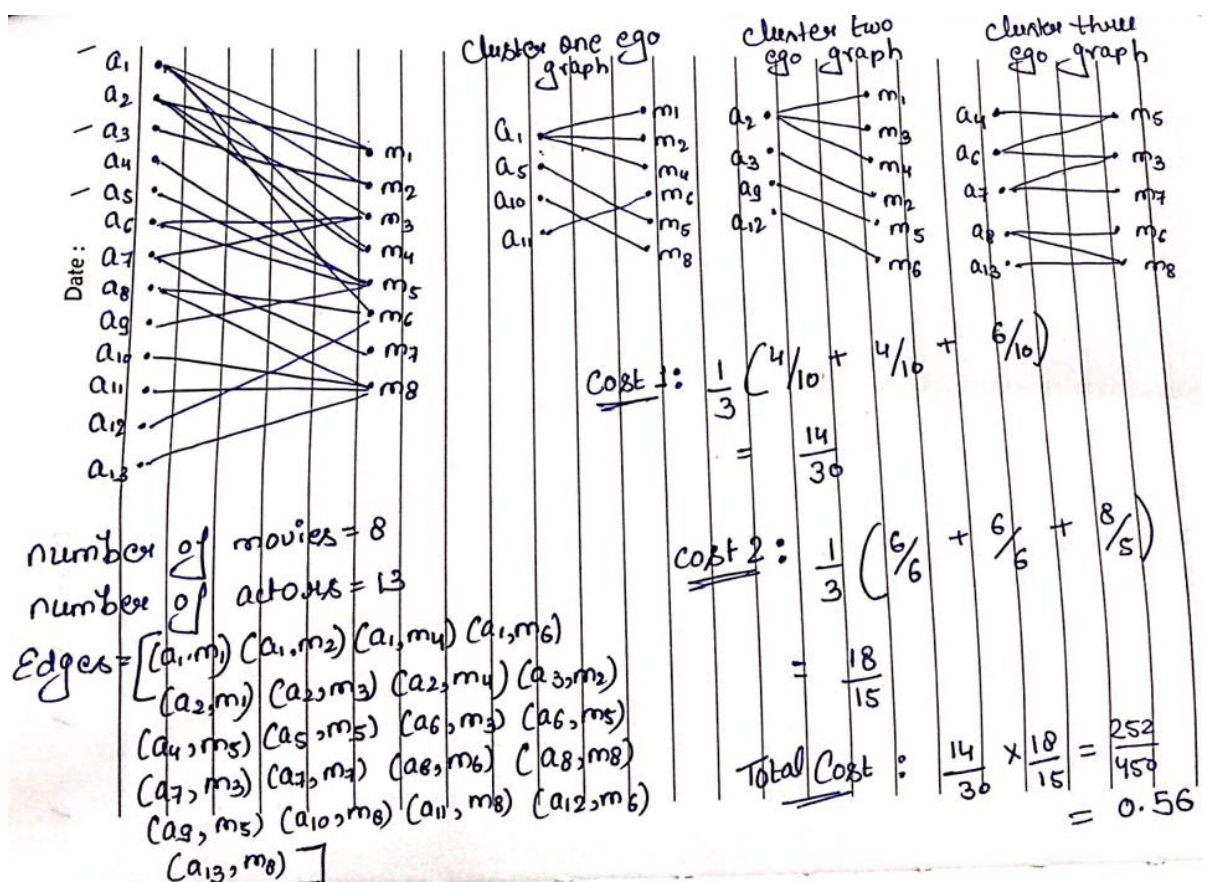
**Every Grader function has to return True.**

**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 =
   $$\frac{1}{N} \sum\nolimits_{\text{each cluster i}} \frac{\text{(number of nodes in the largest connected component in the graph with the actor nodes an}}{\text{(total number of nodes in that cluster i)}}$$
   where N= number of clusters
   (Write your code in def cost1())
5. Cost2 =
   $$\frac{1}{N} \sum\nolimits_{\text{each cluster i}} \frac{\text{(sum of degress of actor nodes in the graph with the actor nodes and its movie neighbours}}{\text{(number of unique movie nodes in the graph with the actor nodes and its movie neighbour}}$$
   where N= 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 =
   $$\frac{1}{N} \sum_{\text{each cluster i}} \frac{\text{(number of nodes in the largest connected component in the graph with the movie nodes a}}{\text{(total number of nodes in that cluster i)}}$$
   where N= number of clusters
   (Write your code in def cost1())

4. Cost2 =
   $$\frac{1}{N} \sum_{\text{each cluster i}} \frac{\text{(sum of degress of movie nodes in the graph with the movie nodes and its actor neighbour}}{\text{(number of unique actor nodes in the graph with the movie nodes and its actor neighbours}}$$
   where N= number of clusters
   (Write your code in def cost2())

**NOTE**: For task1 cost1*cost2 Value should be less than 15 for value n_cluster=3 then it should gradually decrease as n_cluster increases, for task2 value should less than 5 then it should decrease

**Algorithm for actor nodes**

```
        for number_of_clusters in [3, 5, 10, 30, 50, 100, 200,
    500]:
            algo = clustering_algorith(clusters=number_of_clust
    ers)
            # 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 node
    s (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(graph
    2)+cost1(graph3) # here we are doing summation
                cost2=cost2(graph1)+cost2(graph2)+cost2(graph3)
            computer the metric Cost = Cost1*Cost2
        return number_of_clusters which have maximum Cost
```

In [4]:
```python
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
# import stellargraph as sg
# you need to have tensorflow
from stellargraph.data import UniformRandomMetaPathWalk
from stellargraph import StellarGraph
```

In [5]:
```python
data=pd.read_csv('movie_actor_network.csv', index_col=False, names=
```

In [6]:
```python
data.head(5)
```

Out[6]:

| | movie | actor |
|---|---|---|
| **0** | m1 | a1 |
| **1** | m2 | a1 |
| **2** | m2 | a2 |
| **3** | m3 | a1 |
| **4** | m3 | a3 |

In [7]:
```python
edges = [tuple(x) for x in data.values.tolist()]
```

In [8]:
```python
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 [9]:
```python
A = list(nx.connected_component_subgraphs(B))[0]
```

In [10]:
```python
print("number of nodes", A.number_of_nodes())
print("number of edges", A.number_of_edges())
```

```
number of nodes 4703
number of edges 9650
```

In [11]:
```python
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, pos=pos, with_labels=True)
plt.show()
```



In [12]:
```python
movies = []
actors = []
for i in A.nodes():
    if 'm' in i:
        movies.append(i)
    if 'a' in i:
        actors.append(i)
print('number of movies ', len(movies))
print('number of actors ', len(actors))
```

```
number of movies  1292
number of actors  3411
```

In [13]:
```python
# 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
              n=1,         # number of random walks per root node
              metapaths=metapaths
              )

print("Number of random walks: {}".format(len(walks)))
```

```
Number of random walks: 4703
```

In [14]:
```python
from gensim.models import Word2Vec
model = Word2Vec(walks, size=128, window=5)
```

In [15]:
```python
model.wv.vectors.shape  # 128-dimensional vector for each node in t
```

Out[15]:
```
(4703, 128)
```

In [16]:
```python
# Retrieve node embeddings and corresponding subjects
node_ids = model.wv.index2word  # list of node IDs
node_embeddings = model.wv.vectors  # numpy.ndarray of size number
node_targets = [ A.node[node_id]['label'] for node_id in node_ids]
```

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

```
['a973', 'a967', 'a964', 'a1731', 'a969', 'a970', 'a1028', 'a1057', 'a965', 'a1003', 'm1094', 'a966', 'm67', 'a988', 'm1111']
```

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

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

```python
In [18]: def data_split(node_ids,node_targets,node_embeddings):
             '''In this function, we will split the node embeddings into act
             actor_nodes,movie_nodes=[],[]
             actor_embeddings,movie_embeddings=[],[]
             # split the node_embeddings into actor_embeddings,movie_embeddi
             # By using node_ids and node_targets, we can extract actor_node
             for i in range(0,len(node_ids)):
               if node_targets[i]=='movie':
                 movie_embeddings.append(node_embeddings[i])
                 movie_nodes.append(node_ids[i])
               else:
                 actor_embeddings.append(node_embeddings[i])
                 actor_nodes.append(node_ids[i])

             return actor_nodes,movie_nodes,actor_embeddings,movie_embedding
```

```python
In [19]: actor_nodes,movie_nodes,actor_embeddings,movie_embeddings = data_sp
```

Grader function - 1

```python
In [20]: def grader_actors(data):
             assert(len(data)==3411)
             return True
         grader_actors(actor_nodes)
```

Out[20]: True

Grader function - 2

```python
In [21]: def grader_movies(data):
             assert(len(data)==1292)
             return True
         grader_movies(movie_nodes)
```

Out[21]: True

Calculating cost1

Cost1 =

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

where N= number of clusters
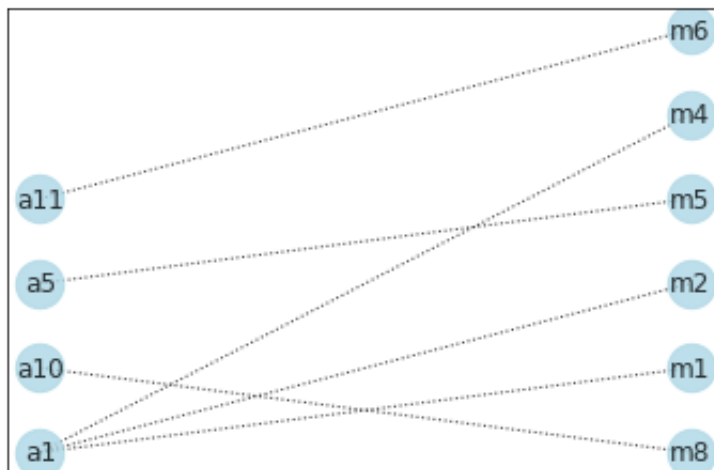
```python
In [22]: def cost1(graph,number_of_clusters):
             '''In this function, we will calculate cost1'''

             connected = nx.connected_components(graph)
             max_cc = max(connected, key=len)
             number_of_nodes_in_cc = len(max_cc)
             total_nodes = graph.number_of_nodes()

             cost1= (1/number_of_clusters) * (number_of_nodes_in_cc / total_

             return cost1
```

```python
In [23]: import networkx as nx
         from networkx.algorithms import bipartite
         graded_graph= nx.Graph()
         graded_graph.add_nodes_from(['a1','a5','a10','a11'], bipartite=0) #
         graded_graph.add_nodes_from(['m1','m2','m4','m6','m5','m8'], bipart
         graded_graph.add_edges_from([('a1','m1'),('a1','m2'),('a1','m4'),('
         l={'a1','a5','a10','a11'};r={'m1','m2','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, pos=pos, with_labels=True,node_color
```



Grader function - 3

```python
In [24]: graded_cost1=cost1(graded_graph,3)
         def grader_cost1(data):
             assert(data==((1/3)*(4/10))) # 1/3 is number of clusters
             return True
         graded_cost1(graded_cost1)
```

Out[24]: True

Calculating cost2

Cost2 =

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

where N= number of clusters

```
In [25]: def cost2(graph,number_of_clusters):
             '''In this function, we will calculate cost1'''

             actors = []
             movies = []
             sum_of_degrees = unique_movies = 0

             for node in graph.nodes():
               if 'a' in node:
                 actors.append(node)
               else:
                 movies.append(node)

             unique_movies = len(movies)
             for a in actors:
               sum_of_degrees += graph.degree(a)

             cost2= (1/number_of_clusters) * (sum_of_degrees/unique_movies)

             return cost2
```

Grader function - 4

```
In [26]: 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[26]: True

Grouping similar actors

In [27]:
```python
from sklearn.cluster import KMeans

number_of_clusters = [3, 5, 10, 30, 50, 100, 200, 500]
Cost={}

for i in number_of_clusters:
  model_k = KMeans(n_clusters=i, random_state=0)
  model_k.fit(actor_embeddings)
  actor_labels = model_k.labels_
  unique_clusters = np.unique(actor_labels)
  dict_of_actor_nodes = dict(zip(actor_nodes, actor_labels))
  list_of_clusters = []
  for n in unique_clusters:
    clusters = []
    for node, cluster in dict_of_actor_nodes.items():
      if cluster == n:
        clusters.append(node)
    list_of_clusters.append(clusters)

  Cost1 = 0
  Cost2 = 0

  for cluster in list_of_clusters:
    G = nx.Graph()
    for node in cluster:
      subgraph = nx.ego_graph(B, node)
      G.add_nodes_from(subgraph.nodes())
      G.add_edges_from(subgraph.edges())

    Cost1 += cost1(G, len(list_of_clusters))
    Cost2 += cost2(G, len(list_of_clusters))

  Cost[i] = Cost1*Cost2
```

In [28]:
```python
Cost
```

Out[28]:
```
{3: 3.7187652319659743,
 5: 3.0106081375880396,
 10: 2.2797206179557974,
 30: 1.7561388855624738,
 50: 1.5005699732498796,
 100: 1.3886519805944508,
 200: 1.6851566363575394,
 500: 1.825831568821166}
```

In [39]:
```python
model = KMeans(n_clusters=3)
model.fit(actor_embeddings)

predict = model.predict(actor_embeddings)
```
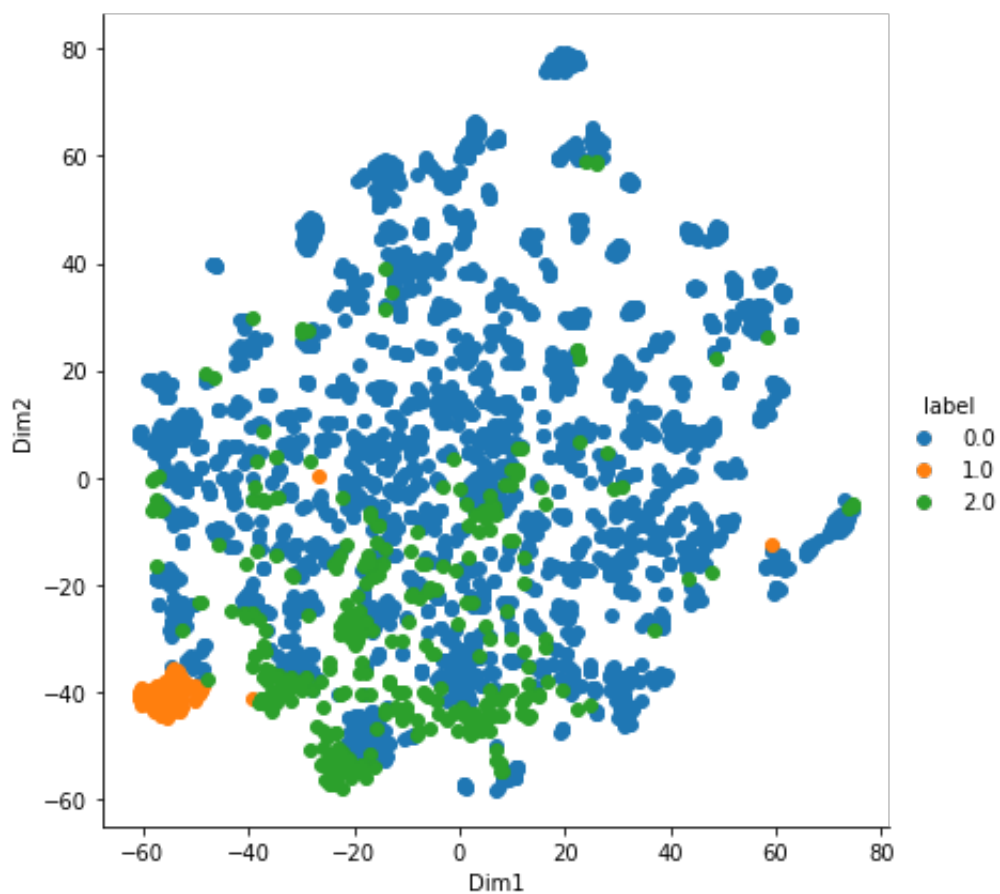
```
In [30]: from sklearn.manifold import TSNE
         transform = TSNE #PCA

         trans = transform(n_components=2)
         twoD_data = trans.fit_transform(actor_embeddings)
```

```
In [31]: req_data = np.vstack((twoD_data.T, predict.T))
         final_data = pd.DataFrame(req_data.T, columns=("Dim1", "Dim2", "lab
```

Displaying similar actor clusters

```
In [32]: import seaborn as sns
         sns.FacetGrid(final_data, hue="label", size=6).map(plt.scatter, "Di
         plt.show()
```



Grouping similar movies

In [33]:
```python
from sklearn.cluster import KMeans

Cost={}

for i in number_of_clusters:
  model_k = KMeans(n_clusters=i, random_state=0)
  model_k.fit(movie_embeddings)
  movie_labels = model_k.labels_
  unique_clusters = np.unique(movie_labels)
  dict_of_movie_nodes = dict(zip(movie_nodes, movie_labels))
  list_of_clusters = []
  for n in unique_clusters:
    clusters = []
    for node, cluster in dict_of_movie_nodes.items():
      if cluster == n:
        clusters.append(node)
    list_of_clusters.append(clusters)

  Cost1 = 0
  Cost2 = 0

  for cluster in list_of_clusters:
    G = nx.Graph()
    for node in cluster:
      subgraph = nx.ego_graph(B, node)
      G.add_nodes_from(subgraph.nodes())
      G.add_edges_from(subgraph.edges())

    Cost1 += cost1(G, len(list_of_clusters))
    Cost2 += cost2(G, len(list_of_clusters))

  Cost[i] = Cost1*Cost2
```

In [34]:
```python
Cost
```

Out[34]:
```
{3: 8.428757308074937,
 5: 8.768042171312553,
 10: 9.2982585656088,
 30: 12.581764604004045,
 50: 12.419826998888821,
 100: 13.740654302675953,
 200: 12.589023573602638,
 500: 10.318293861592426}
```

In [40]:
```python
model = KMeans(n_clusters=100)
model.fit(movie_embeddings)

predict = model.predict(movie_embeddings)
```

In [36]:
```python
from sklearn.manifold import TSNE
transform = TSNE #PCA

trans = transform(n_components=2)
twoD_data = trans.fit_transform(movie_embeddings)
```

In [37]:
```python
req_data = []
final_data = []
req_data = np.vstack((twoD_data.T, predict.T))
final_data = pd.DataFrame(req_data.T, columns=("Dim1", "Dim2", "lab
```

Displaying similar movie clusters

In [38]:
```python
import seaborn as sns
sns.FacetGrid(final_data, hue="label", size=6).map(plt.scatter, "Di
plt.show()
```