

## ▼ 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.

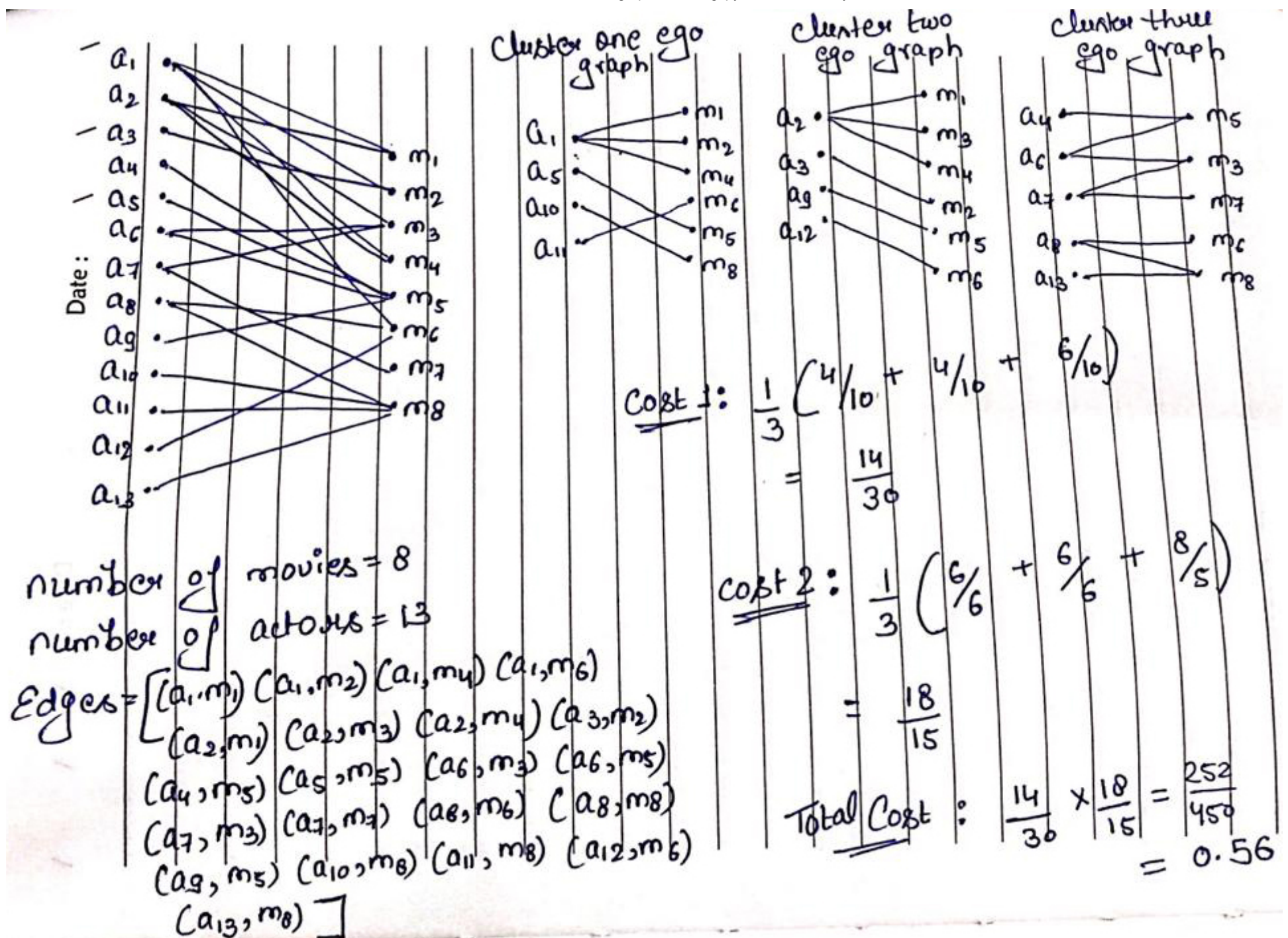
Every Grader function has to return True.

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

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>
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= 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= 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  $Cost1 * Cost2$

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

where N= number of clusters

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

$$Cost2 = \frac{1}{N} \sum_{\text{each cluster } i} \frac{(\text{sum of degrees of movie nodes in the graph with the movie nodes and its actor neighbours in cluster } i)}{(\text{number of unique actor nodes in the graph with the movie nodes and its actor neighbours in cluster } i)}$$

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(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
  
```

## Importing libraries

```
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
```

## Importing dataset

```
1 from google.colab import files
2 uploaded = files.upload()
```

 **Choose Files** movie\_actor\_network.csv

- **movie\_actor\_network.csv**(application/vnd.ms-excel) - 114060 bytes, last modified: 5/3/2020 - 100% done  
Saving movie\_actor\_network.csv to movie\_actor\_network.csv


```
1 data=pd.read_csv('movie_actor_network.csv', index_col=False, names=['movie','actor'])
```

```
1 edges = [tuple(x) for x in data.values.tolist()]
```

```
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')
```

```
1 A = list(nx.connected_component_subgraphs(B))[0]
```

```
1 print("number of nodes", A.number_of_nodes())
2 print("number of edges", A.number_of_edges())
```

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

```
1 l, r = nx.bipartite.sets(A)
2 pos = {}
3
4 pos.update((node, (1, index)) for index, node in enumerate(l))
5 pos.update((node, (2, index)) for index, node in enumerate(r))
6
7 nx.draw(A, pos=pos, with_labels=True)
8 plt.show()
```



```
1 movies = []
2 actors = []
3 for i in A.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

```

```

1
2 # Create the random walker
3 rw = UniformRandomMetaPathWalk(StellarGraph(A))
4
5 # specify the metapath schemas as a list of lists of node types.
6 metapaths = [
7     ["movie", "actor", "movie"],
8     ["actor", "movie", "actor"]
9 ]
10
11 walks = rw.run(nodes=list(A.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

```

```

1 from gensim.models import Word2Vec
2 model = Word2Vec(walks, size=128, window=5)

```

```

1 model.wv.vectors.shape # 128-dimensional vector for each node in the graph

```

```

↳ (4703, 128)

```

```

1 # Retrieve node embeddings and corresponding subjects
2 node_ids = model.wv.index2word # list of node IDs
3 node_embeddings = model.wv.vectors # numpy.ndarray of size number of nodes times embeddings dimensionality
4 node_targets = [ A.node[node_id]['label'] for node_id in node_ids]

```

```

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

```

```

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

```

```

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

```

```

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

```

```

1 def data_split(node_ids,node_targets,node_embeddings):
2     '''In this function, we will split the node embeddings into actor_embeddings , movie_embeddings '''
3     actor_nodes,movie_nodes=[],[]
4     actor_embeddings,movie_embeddings=[],[]
5     for i in range(len(node_ids)):
6         if 'm' in node_targets[i]:
7             movie_nodes.append(node_ids[i])
8             movie_embeddings.append(node_embeddings[i])
9         else:
10            actor_nodes.append(node_ids[i])
11            actor_embeddings.append(node_embeddings[i])
12
13
14     return actor_nodes,movie_nodes,actor_embeddings,movie_embeddings
15

```

```

1 actor_nodes,movie_nodes,actor_embeddings,movie_embeddings=data_split(node_ids,node_targets,node_embeddings)

```

## Graded function - 1

```

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

```



 True

Graded function - 2

```
1 def grader_movies(data):
2     assert(len(data)==1292)
3     return True
4 grader_movies(movie_nodes)
```

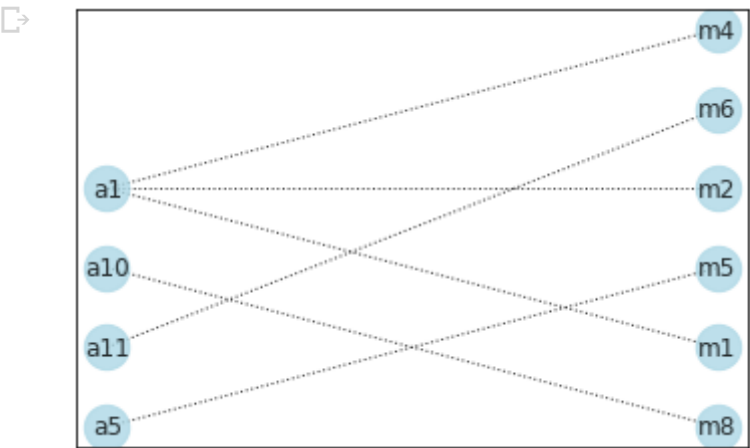
 True

Calculating cost1

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

```
1 def cost_1(graph,cluster):
2     Gc = max(nx.connected_component_subgraphs(graph), key=len)
3     connected=Gc.number_of_nodes()
4     total_nodes=graph.number_of_nodes()
5     value=((1/cluster)*(connected/total_nodes))
6     return value

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 node attribute "bipartite"
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'),('a5','m5'),('a10','m8')])
7 l={'a1','a5','a10','a11'};r={'m1','m2','m4','m6','m5','m8'}
8 pos = {}
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',alpha=0.8,style='dotted',node_size=500)
```



Graded function - 3

```
1 graded_cost1=cost_1(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)
```

 True

Calculating cost2

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

```
1 def cost_2(graph,cluster):
2     degree=graph.degree()
3     degree_value=0 # to store sum of all  degree of actor nodes
4     mov_nodes=0 #to store sum of all unique actor nodes
5     for i in degree:
6         if "a" in i[0]:
7             degree_value=degree_value+i[1]
8     else:
```

```

-      -----
9      mov_nodes=mov_nodes+1
10     value=((1/cluster)*(degree_value/mov_nodes))
11     return value
12

```

## Graded function - 4

```

1 graded_cost2=cost_2(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)

```

```

↳ True

```

## Grouping similar actors

```

1 from sklearn.cluster import KMeans

1 cost_value={}
2 for clu in [3,5,10,50,200,300]:
3     cost1=0
4     cost2=0
5     label=[]
6     algo = KMeans(n_clusters=clu)
7     algo.fit(actor_embeddings) # FITTING WITH KMEANS ALGO
8     label=algo.labels_
9     for i in range(clu):
10         G1=nx.Graph() # drawing graph for each cluster
11         label_divi=[]
12         k=[index for index, value in enumerate(label) if value == i]# accesing each cluster based on labels_
13         label_divi=[ actor_nodes[l] for l in k]
14         for node in label_divi:
15             sub_graph1=nx.ego_graph(B,node)
16             G1.add_nodes_from(sub_graph1.nodes) # adding nodes
17             G1.add_edges_from(sub_graph1.edges()) # adding edges
18             cst1=cost_1(G1,clu)
19             cost1=cost1+cst1 # calculating cost functions
20             cst2=cost_2(G1,clu)
21             cost2=cost2+cst2
22         value=cost1*cost2
23         cost_value[clu]=value

```

## Fitting the clustering algorithm with the opimal number\_of\_clusters

```

1 optimal_cluster_number = max(cost_value, key=cost_value.get)
2 print("optimal cluster number = ",optimal_cluster_number) #https://www.geeksforgeeks.org/python-get-key-with-maximum-value-in-dict

```

```

↳ optimal cluster number = 3

```

```

1 model= KMeans(n_clusters=optimal_cluster_number)
2 model.fit(actor_embeddings)
3 label=model.labels_

1 node_cluster={}
2 for i in range(optimal_cluster_number):
3     k=[index for index, value in enumerate(label) if value == i] # getting the cluster number for each node
4     label_divi=[ actor_nodes[l] for l in k]
5     for node in label_divi:
6         node_cluster[node]=i

```

```

1 node_cluster['a1435']

```

```

↳ 2

```

## Displaying similar actor clusters

```

1 from sklearn.manifold import TSNE
2 transform = TSNE #PCA
3 trans = transform(n_components=2)
4 actor_embeddings_2d = trans.fit_transform(actor_embeddings)

```

```

1 import matplotlib.pyplot as plt

```

```

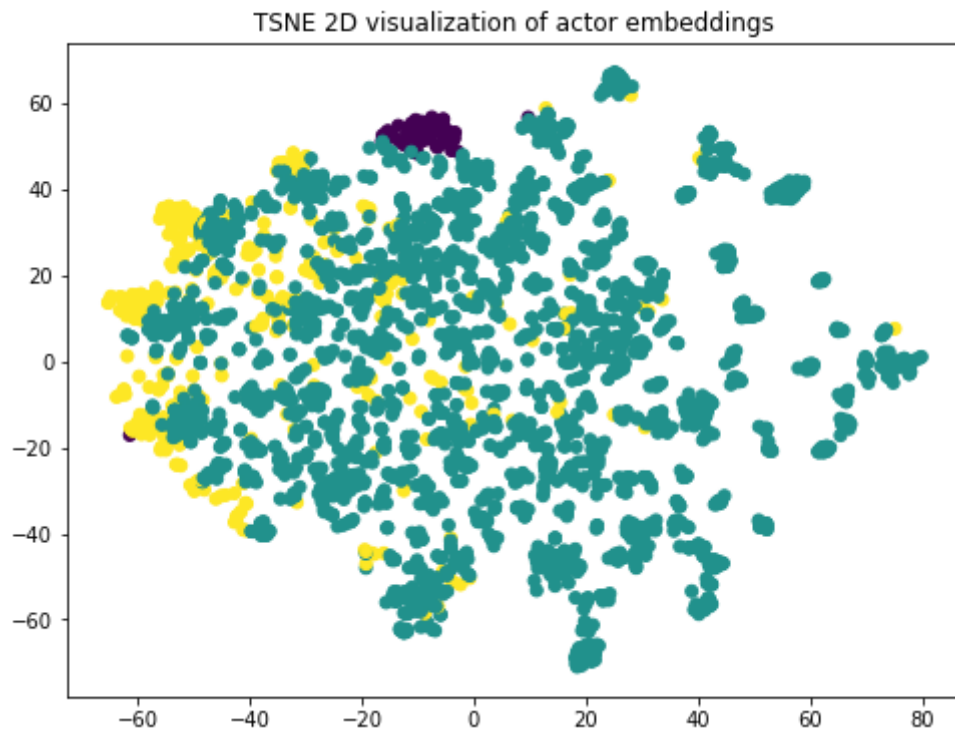
1 import matplotlib.pyplot as plt
2 plt.figure(figsize=(8, 6)) #https://stackoverflow.com/questions/28227340/kmeans-scatter-plot-plot-different-colors-per-cluster
3 plt.scatter(actor_embeddings_2d[:,0], actor_embeddings_2d[:,1], c=model.labels_.astype(float))
4 plt.title('TSNE 2D visualization of actor embeddings')

```

```

↳ Text(0.5, 1.0, 'TSNE 2D visualization of actor embeddings')

```



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

```

1 def cost_2(graph,cluster):
2     degree=graph.degree()
3     degree_value=0 # to store sum of all degree of actor nodes
4     act_nodes=0 #to store sum of all unique actor nodes
5     for i in degree:
6         if "m" in i[0]:
7             degree_value=degree_value+i[1]
8         else:
9             act_nodes=act_nodes+1
10    value=((1/cluster)*(degree_value/act_nodes))
11    return value

```

### Grouping similar movies

```

1 cost_value={}
2 for clu in [3,5,10,50,200,300]:
3     cost1=0
4     cost2=0
5     label=[]
6     algo = KMeans(n_clusters=clu)
7     algo.fit(movie_embeddings) # FITTING WITH KMEANS ALGO
8     label=algo.labels_
9     for i in range(clu):
10        G1=nx.Graph() # drawing graph for each cluster
11        label_divi=[]
12        k=[index for index, value in enumerate(label) if value == i]# accesing each cluster based on labels_
13        label_divi=[ movie_nodes[l] for l in k]
14        for node in label_divi:
15            sub_graph1=nx.ego_graph(B,node)
16            G1.add_nodes_from(sub_graph1.nodes) # adding nodes
17            G1.add_edges_from(sub_graph1.edges()) # adding edges
18        cst1=cost_1(G1,clu)
19        cost1=cost1+cst1 # calculating cost functions
20        cst2=cost_2(G1,clu)
21        cost2=cost2+cst2
22    value=cost1*cost2
23    cost_value[clu]=value

```

### Fitting the clustering algorithm with the opimal number\_of\_clusters

```

1 optimal_cluster_number = max(cost_value, key=cost_value.get)
2 print("optimal cluster number =",optimal_cluster_number) #https://www.geeksforgeeks.org/python-get-key-with-maximum-value-in-dicti

```

```

↳ optimal cluster number = 3

```

```

1 model= KMeans(n_clusters=optimal_cluster_number)

```

```

2 model.fit(movie_embeddings)
3 label=model.labels_

1 node_cluster={}
2 for i in range(optimal_cluster_number):
3     k=[index for index, value in enumerate(label) if value == i]
4     label_divi=[ movie_nodes[l] for l in k]
5     for node in label_divi:
6         node_cluster[node]=i

```

```

1 node_cluster['m858']

```

1

Displaying similar movie clusters

```

1 from sklearn.manifold import TSNE
2 transform = TSNE #PCA
3 trans = transform(n_components=2)
4 movie_embeddings_2d = trans.fit_transform(movie_embeddings)

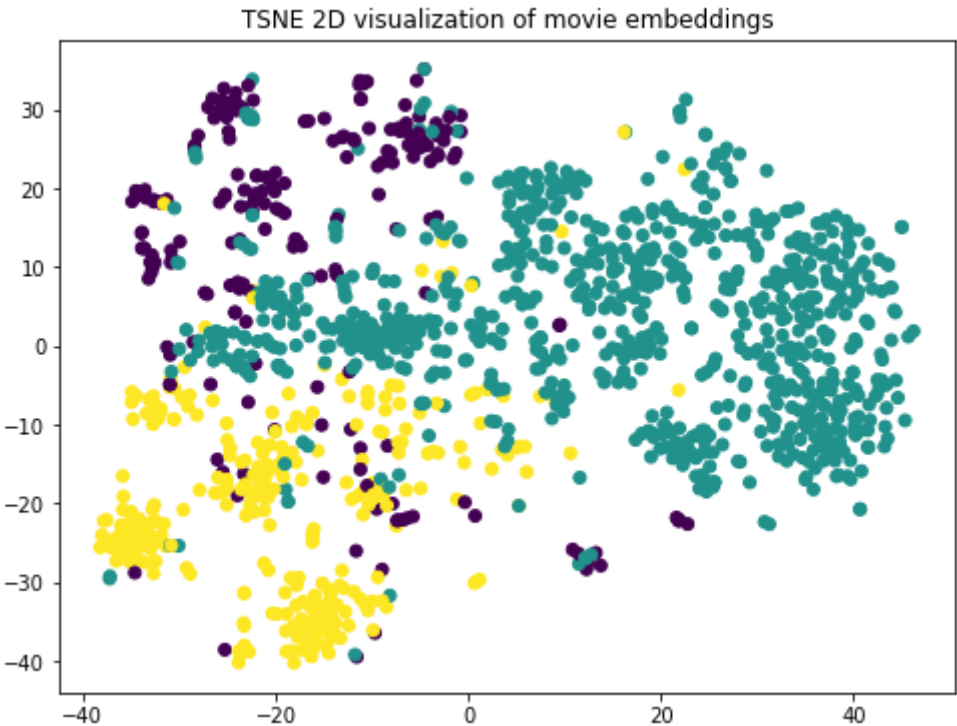
1 import matplotlib.pyplot as plt
2 # Visualize it:
3 plt.figure(figsize=(8, 6))
4 plt.scatter(movie_embeddings_2d[:,0], movie_embeddings_2d[:,1], c=model.labels_.astype(float))
5 plt.title('TSNE 2D visualization of movie embeddings')

```

```

Text(0.5, 1.0, 'TSNE 2D visualization of movie embeddings')

```



1