

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.

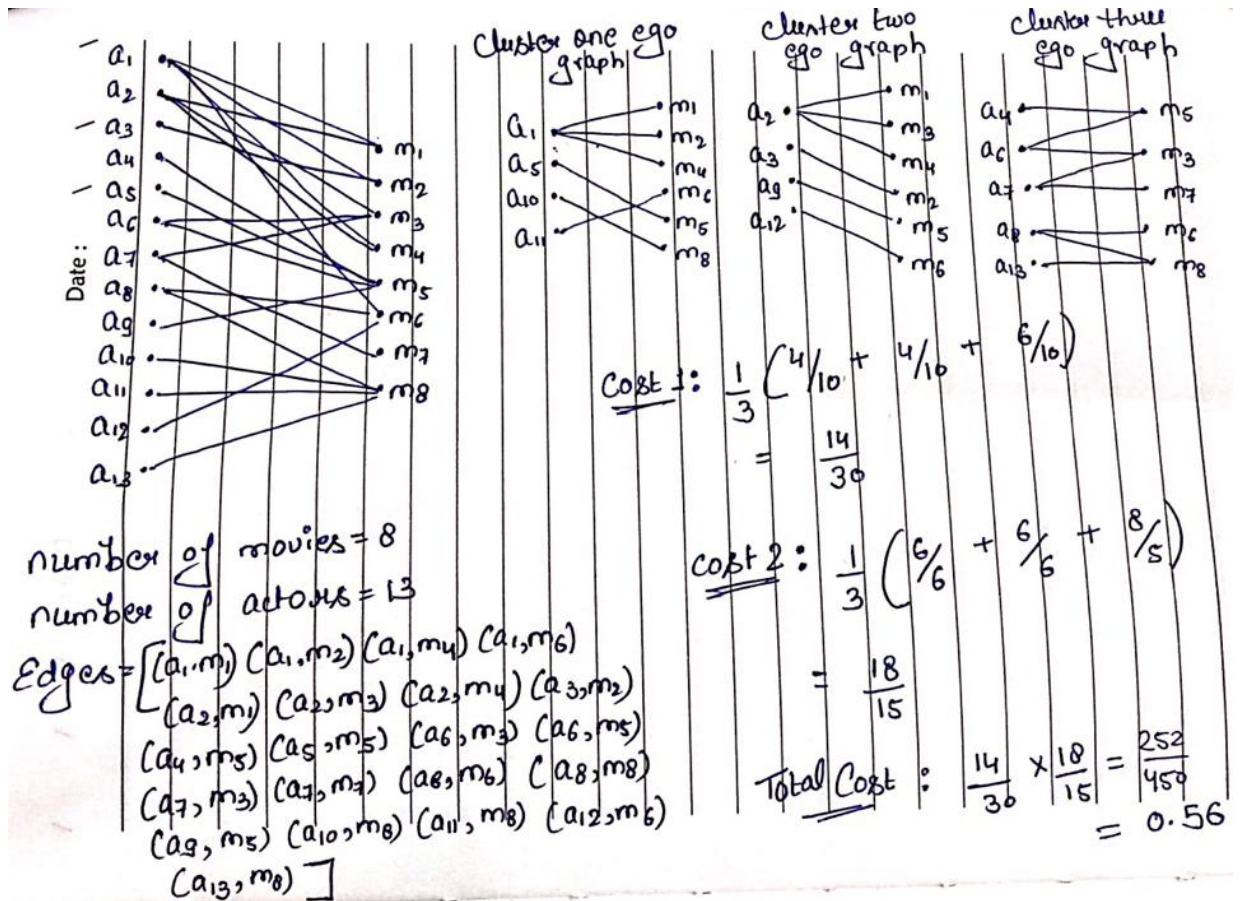
Please check [clustering assignment helper functions](https://drive.google.com/file/d/1V29KhKo3YnckMX32treEgdtH5r90DljU/view?usp=sharing)

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

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

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 act
or 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
h3) # 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 [1]: !pip install networkx==2.3

```

Requirement already satisfied: networkx==2.3 in c:\users\disha\anaconda3\lib\site-packages (2.3)
Requirement already satisfied: decorator>=4.3.0 in c:\users\disha\anaconda3\lib\site-packages (from networkx==2.3) (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

```

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

In [3]: 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')

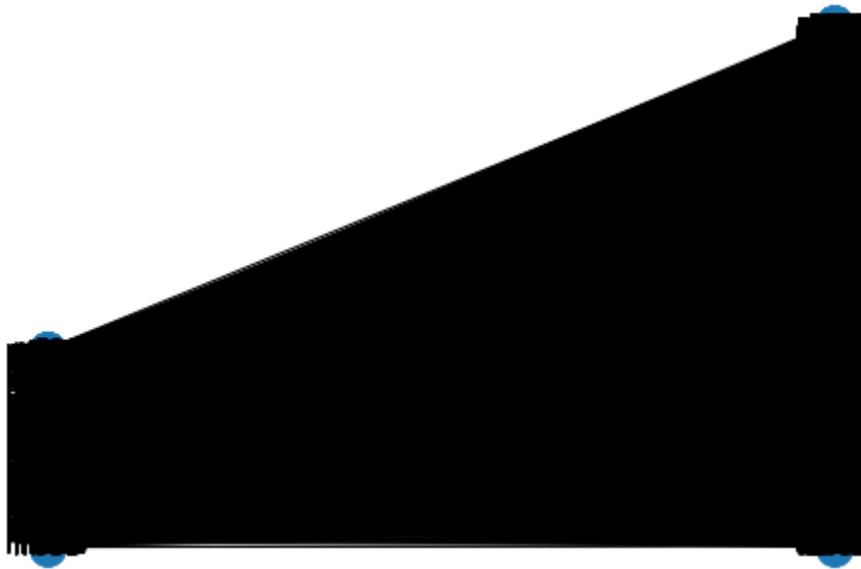
```

In [5]: A = list(nx.connected_component_subgraphs(B))[0]

```
In [6]: 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 [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, pos=pos, with_labels=True)  
        plt.show()
```



```
In [8]: 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 [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
               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 [10]:

```

from gensim.models import Word2Vec
model = Word2Vec(walks, 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.wv.index2word # List of node IDs
node_embeddings = model.wv.vectors # numpy.ndarray of size number of nodes times
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']

```

```
In [13]: import re
def data_split(node_ids,node_targets,node_embeddings):
    '''In this function, we will split the node embeddings into actor_embeddings
    actor_nodes,movie_nodes=[],[]
    actor_embeddings,movie_embeddings=[],[]
    # split the node_embeddings into actor_embeddings,movie_embeddings based on r
    # By using node_embedding and node_targets, we can extract actor_embedding ar
    # By using node_ids and node_targets, we can extract actor_nodes and movie no

    i = 0
    while i < len(node_ids):
        if re.search('^a',node_ids[i]):
            actor_nodes.append(node_ids[i])
            actor_embeddings.append(node_embeddings[i])
        else:
            movie_nodes.append(node_ids[i])
            movie_embeddings.append(node_embeddings[i])
        i = i + 1

    return actor_nodes,movie_nodes,actor_embeddings,movie_embeddings
```

```
In [14]: actor_nodes,movie_nodes,actor_embeddings,movie_embeddings=data_split(node_ids,noc
```

Grader function - 1

```
In [15]: def grader_actors(data):
    assert(len(data)==3411)
    return True
grader_actors(actor_nodes)
```

Out[15]: True

Grader function - 2

```
In [16]: def grader_movies(data):
    assert(len(data)==1292)
    return True
grader_movies(movie_nodes)
```

Out[16]: True

Calculating cost1

Cost1 =

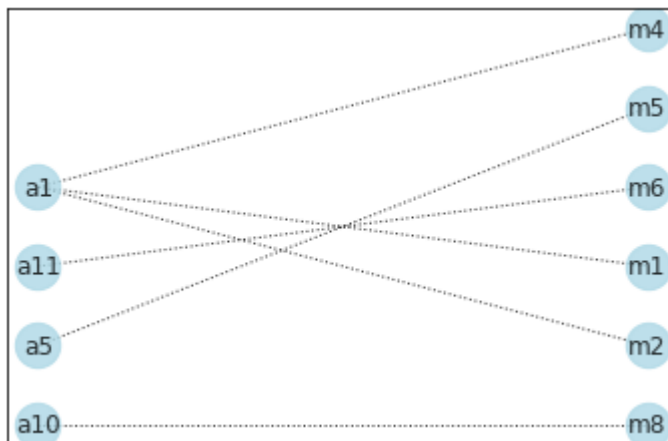
$$\frac{1}{N} \sum_{\text{each cluster } i} \frac{(\text{number of nodes in the largest connected component in the graph with the actor nodes and its movie neighbors})}{(\text{total number of nodes in that cluster } i)}$$

where N= number of clusters

```
In [55]: def cost1(graph,number_of_clusters):
'''In this function, we will calculate cost1'''
nodes = graph.nodes
a = nx.connected_components(graph)
maximum = max([len(i) for i in a])
cost1 = (1/number_of_clusters)*((maximum)/graph.number_of_nodes())

return cost1
```

```
In [56]: 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 nodes
graded_graph.add_nodes_from(['m1','m2','m4','m6','m5','m8'], bipartite=1)
graded_graph.add_edges_from([('a1','m1'),('a1','m2'),('a1','m4'),('a11','m6'),('a11','m5'),('a11','m8'),('a5','m4'),('a5','m6'),('a5','m8'),('a10','m1'),('a10','m2'),('a10','m4'),('a10','m5'),('a10','m8')])
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='lightblue',a
```



Grader function - 3

```
In [57]: 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[57]: True

Calculating cost2

Cost2 =

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

where N= number of clusters

```
In [58]: def cost2(graph,number_of_clusters):
'''In this function, we will calculate cost1'''
nodes = graph.nodes
top_nodes = {n for n, d in list(graph.nodes(data=True)) if d["bipartite"] == 1}
bottom_nodes = nodes - top_nodes
degree = 0
for i in top_nodes:
    degree = degree + graph.degree[i]
cost2= (1/number_of_clusters) * (degree/len(bottom_nodes))

return cost2
```

Grader function - 4

```
In [59]: 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[59]: True

Grouping similar actors

```
In [60]: K = [3, 5, 10, 30, 50, 100, 200, 500]
cost = []
c = dict()
for k in [3, 5, 10, 30, 50, 100, 200, 500]:
    kmeans = KMeans(n_clusters=k, random_state=0).fit(actor_embeddings)
    labels = kmeans.labels_
    d = dict()
    for i in range(len(actor_nodes)):
        d[actor_nodes[i]] = labels[i]
    clusters = []
    for i in range(k):
        e = [k1 for k1,v in d.items() if v == i]
        clusters.append(e)
    graphs = []
    for i in range(k):
        g = nx.Graph()
        for n in clusters[i]:
            graph = nx.ego_graph(B,n)
            g.add_nodes_from(graph.nodes,bipartite = 1)
            g.add_edges_from(graph.edges)
            g.add_nodes_from(clusters[i],bipartite = 0)
            graphs.append(g)
    c1,c2 = 0,0
    for i in range(k):
        c1 = c1 + (cost1(graphs[i],k))
        c2 = c2 + (cost2(graphs[i],k))
    cost.append((c1)*(c2))
    c[k] = c1*c2
```



```
In [61]: ind = cost.index(max(cost))  
best_k = K[ind]
```

```
In [62]: c
```

```
Out[62]: {3: 3.8863795023705694,  
5: 2.981023298207314,  
10: 2.174437913426483,  
30: 1.6117862118008819,  
50: 1.5685894889997078,  
100: 1.58875720918826,  
200: 1.7012620344490466,  
500: 1.8384617856234164}
```

```
In [63]: best_k
```

```
Out[63]: 3
```

```
In [64]: model = KMeans(n_clusters=best_k, random_state=0).fit(actor_embeddings)  
label = model.labels_
```

```
In [65]: from sklearn.manifold import TSNE  
transform = TSNE  
  
trans = transform(n_components=2)  
node_embeddings_2d = trans.fit_transform(actor_embeddings)
```

```
In [66]: label
```

```
Out[66]: array([1, 1, 1, ..., 2, 2, 2])
```

```
In [ ]:
```

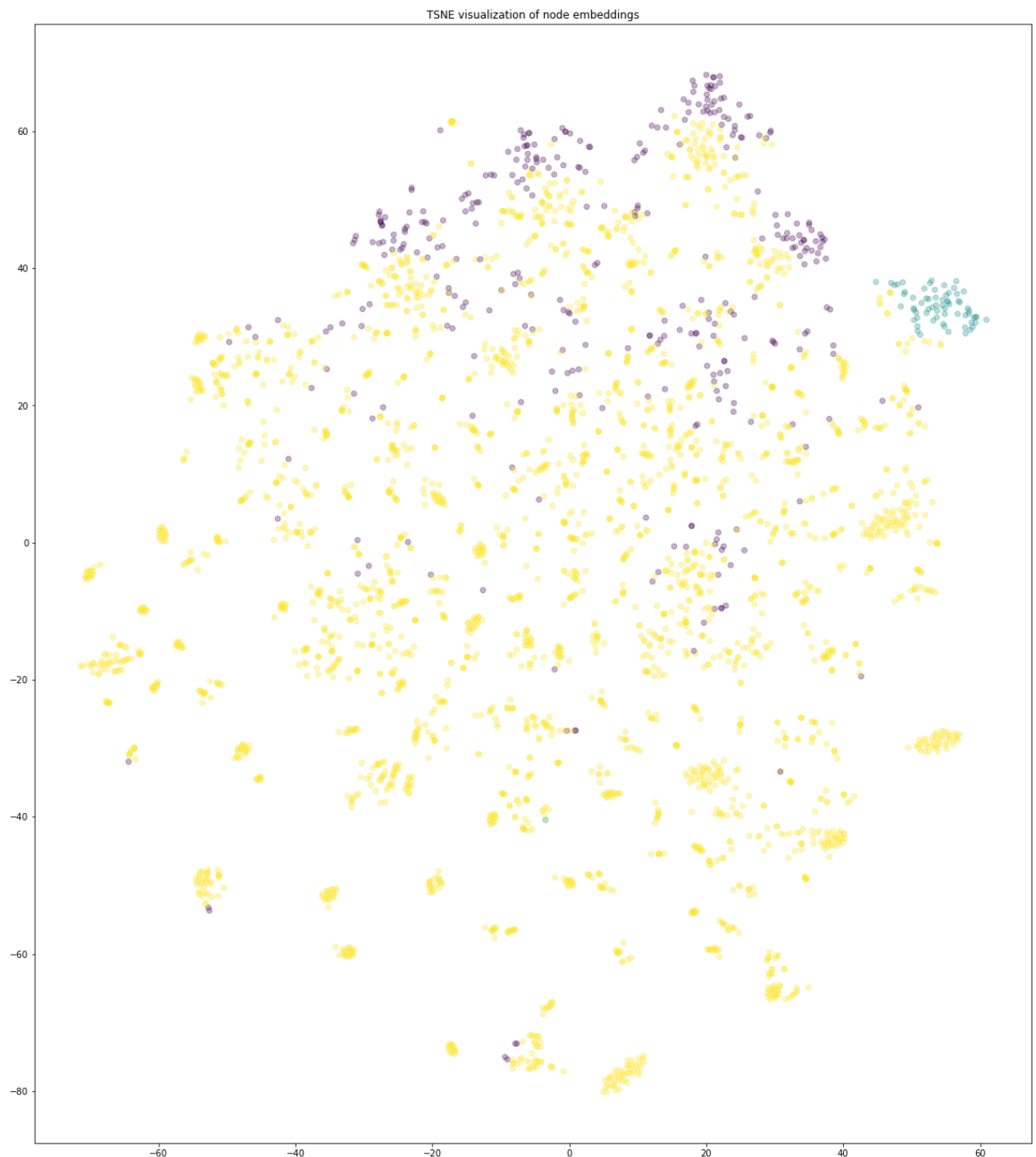
Displaying similar actor clusters

```
In [67]: import numpy as np
# draw the points

label_map = { l: i for i, l in enumerate(np.unique(label))}
node_colours = [ label_map[target] for target in label]

plt.figure(figsize=(20,60))
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()
```



Grouping similar movies

```

In [68]: K = [3, 5, 10, 30, 50, 100, 200, 500]
cost = []
c = dict()
for k in [3, 5, 10, 30, 50, 100, 200, 500]:
    kmeans = KMeans(n_clusters=k, random_state=0).fit(movie_embeddings)
    labels = kmeans.labels_
    d = dict()
    for i in range(len(movie_nodes)):
        d[movie_nodes[i]] = labels[i]
    clusters = []
    for i in range(k):
        e = [k1 for k1,v in d.items() if v == i]
        clusters.append(e)
    graphs = []
    for i in range(k):
        g = nx.Graph()
        for n in clusters[i]:
            graph = nx.ego_graph(B,n)
            g.add_nodes_from(graph.nodes,bipartite = 1)
            g.add_edges_from(graph.edges)
            g.add_nodes_from(clusters[i],bipartite = 0)
        graphs.append(g)
    c1,c2 = 0,0
    for i in range(k):
        c1 = c1 + cost1(graphs[i],k)
        c2 = c2 + cost2(graphs[i],k)
    cost.append(c1*c2)
    c[k] = c1*c2

```

```

In [69]: ind = cost.index(max(cost))
best_k = K[ind]

```

In [70]: best_k

Out[70]: 3

In [71]: c

Out[71]: {3: 2.8352802150560197,
5: 2.5293719381225657,
10: 2.7623660667986893,
30: 2.099404388318737,
50: 1.8586471995041023,
100: 1.5901069039311393,
200: 1.3676762184723799,
500: 1.2116205928020927}

In [72]: `from sklearn.manifold import TSNE`
`transform = TSNE`

`trans = transform(n_components=2)`
`node_embeddings_2d = trans.fit_transform(actor_embeddings)`

In []:

In []:

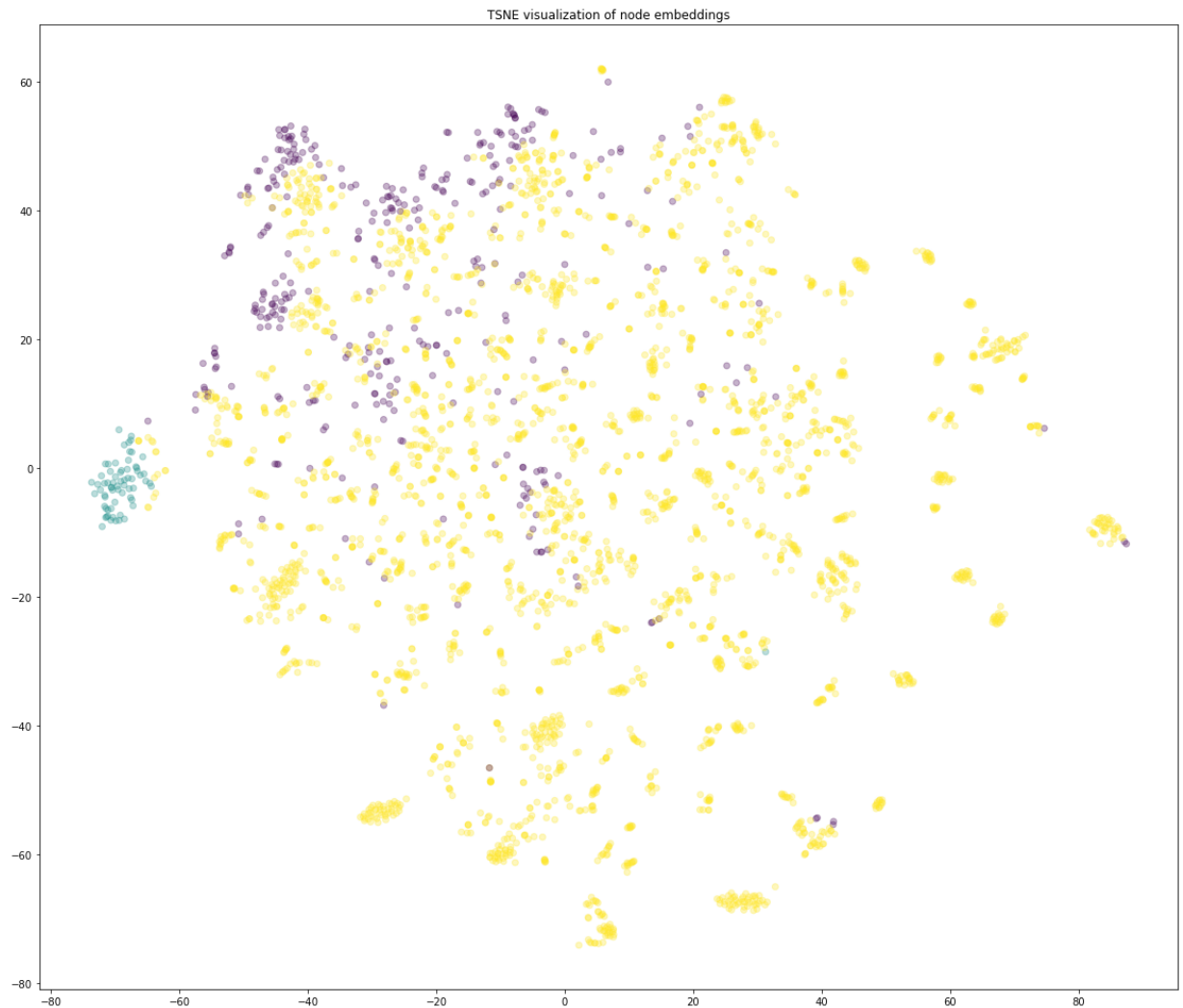
Displaying similar movie clusters

```
In [73]: import numpy as np
# draw the points

label_map = { l: i for i, l in enumerate(np.unique(label))}
node_colours = [ label_map[target] for target in label]

plt.figure(figsize=(20,60))
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 []:

References

reference notebook that was given

<https://networkx.org/documentation/stable/tutorial.html>
(<https://networkx.org/documentation/stable/tutorial.html>)

In []: