

## ▼ 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()`)

```
dict: label_map
```

```
(1 item) {'movie': 0}
```

## ▼ 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.  $Cost1 =$

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

where N= number of clusters

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

5.  $Cost2 =$

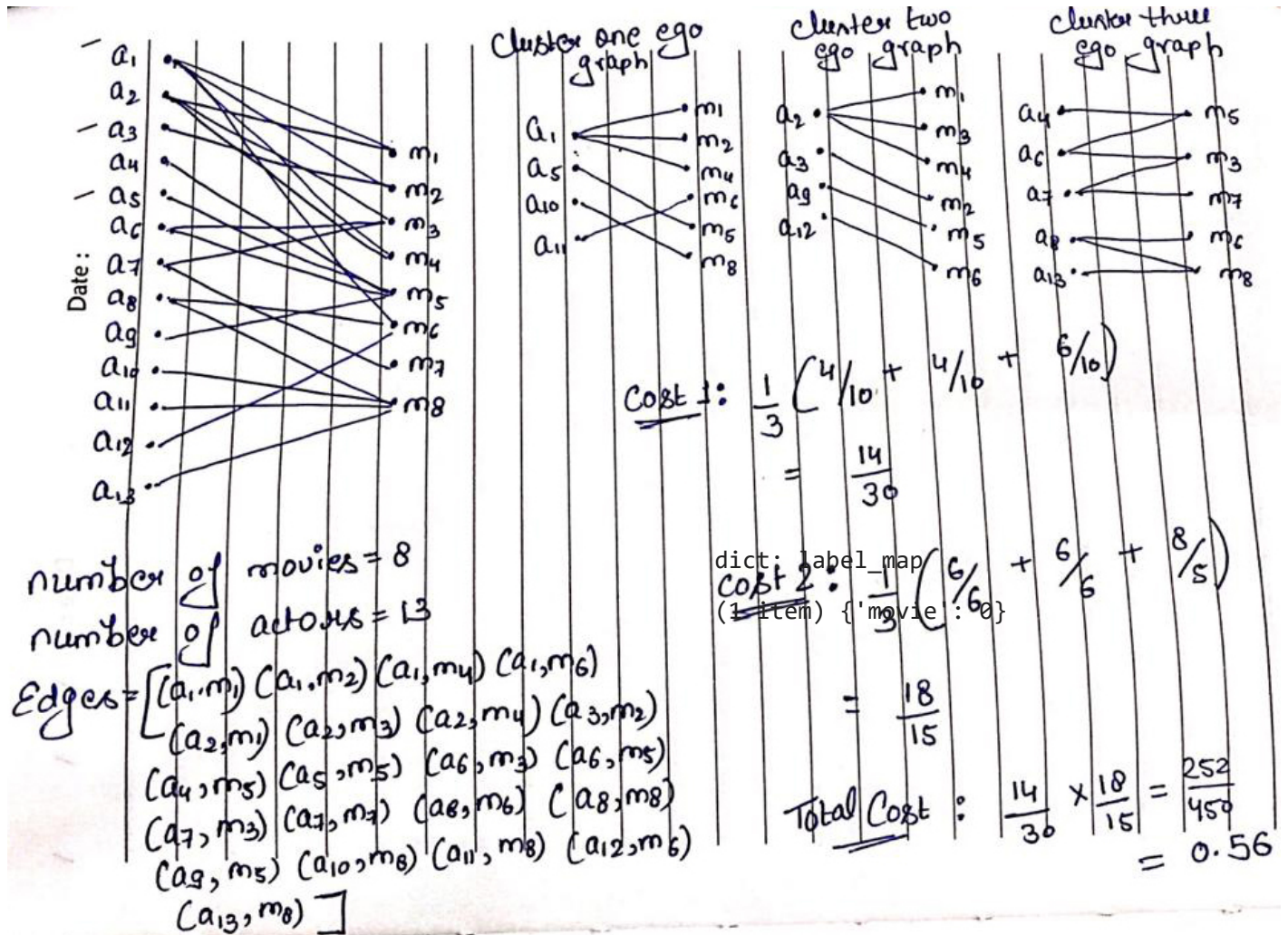
$$\frac{1}{N} \sum_{\text{each cluster } i} \frac{(\text{sum of degress of actor nodes in the graph with the actor nodes and its movie neighbours } i)}{(\text{number of unique movie nodes in the graph with the actor nodes and its movie neighbours } i)}$$

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

- 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 e 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
- Choose the number of clusters for which you have maximum score of  $Cost1 * Cost2$

Cost1 =

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





```
Requirement already satisfied: google-auth<3,>=1.6.3 in /usr/local/lib/python3.7/dist-packages (1.6.3)
Requirement already satisfied: markdown<=2.6.8 in /usr/local/lib/python3.7/dist-packages (2.6.8)
Requirement already satisfied: requests<3,>=2.21.0 in /usr/local/lib/python3.7/dist-packages (2.21.0)
Requirement already satisfied: cachetools<5.0,>=2.0.0 in /usr/local/lib/python3.7/dist-packages (2.0.0)
Requirement already satisfied: pyasn1-modules<=0.2.1 in /usr/local/lib/python3.7/dist-packages (0.2.1)
Requirement already satisfied: rsa<5,>=3.1.4 in /usr/local/lib/python3.7/dist-packages (3.1.4)
Requirement already satisfied: requests-oauthlib<=0.7.0 in /usr/local/lib/python3.7/dist-packages (0.7.0)
Requirement already satisfied: importlib-metadata<=4.4 in /usr/local/lib/python3.7/dist-packages (4.4)
Requirement already satisfied: zipp<=0.5 in /usr/local/lib/python3.7/dist-packages (0.5)
Requirement already satisfied: pyasn1<0.5.0,>=0.4.6 in /usr/local/lib/python3.7/dist-packages (0.4.6)
Requirement already satisfied: certifi<=2017.4.17 in /usr/local/lib/python3.7/dist-packages (2017.4.17)
```

```
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
from stellargraph.data import UniformRandomMetaPathWalk
from stellargraph import StellarGraph
```

dict: label\_map  
(1 item) {'movie': 0}

➤ #Lets get the dataset that has given movie\_actors\_network

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

```
#Defining Edges of the Graph and considering them from tuple
edges = [tuple(x) for x in data.values.tolist()]
```

```
#Lets add Nodes for the Grapg as we already defined edges for them
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')
```

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

```
#Lets print number of nodes and edges got formed from the dataset that has given
print("number of nodes", A.number_of_nodes())
print("number of edges", A.number_of_edges())
```

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

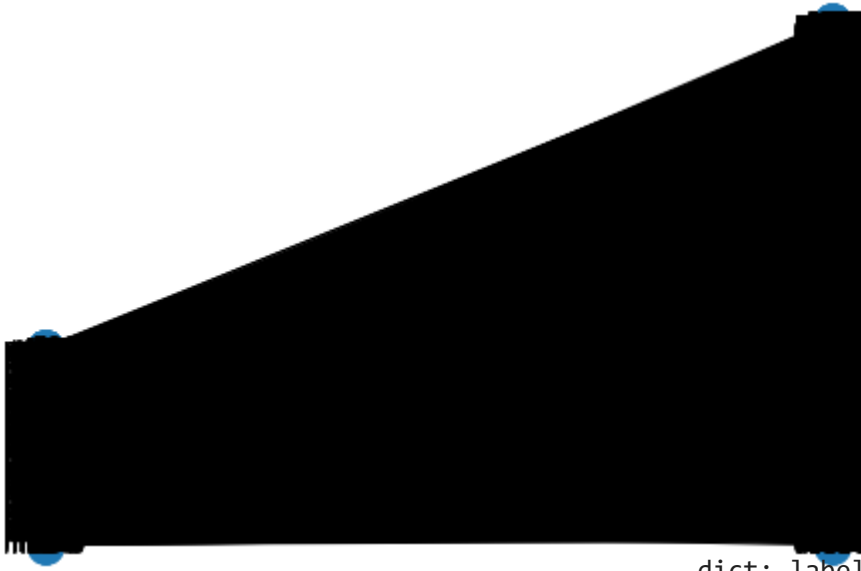
```

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()

```



```

dict: label_map
(1 item) {'movie': 0}

```

```

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

```

```

#https://www.geeksforgeeks.org/random-walk-implementation-python/
# Creating the random walker
#https://stellargraph.readthedocs.io/en/stable/demos/embeddings/metapath2vec-embeddings.html
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, # max length of a random walk
               n=1, # no of random walks per root node
               metapaths=metapaths
               )
```

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

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

```
#lets convert all the words formats in vector forms
from gensim.models import Word2Vec
model = Word2Vec(walks, size=128, window=5)
```

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

(4703, 128)
```

```
# lets 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 embeddings
node_targets = [ A.node[node_id]['label'] for node_id in node_ids]
               dict: label_map
```

```
print(node_ids[:15], end='') (1 item) {'movie': 0}
['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']
```

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

```
['a973', 'a967', 'a964', 'a1731', 'a970', 'a969', 'a1028', 'm1094', 'a1003', 'a965', 'a
```



```
#lets print number of node targets
print("")
print(node_targets[:15],end='')
```

```
['actor', 'actor', 'actor', 'actor', 'actor', 'actor', 'actor', 'actor', 'movie', 'actor', 'acto
```



```
def data_split(node_ids,node_targets,node_embeddings):
    '''In this function, we will split the node embeddings into actor_embeddings , movie_embeddings'''
```

```

actor_nodes, movie_nodes = [], []
actor_embeddings, movie_embeddings = [], []
# splitting the node_embeddings into actor_embeddings and movie_embeddings based on node_i
# using node_embeddings and node_targets, we can even extract actor_embedding and movie e
# By using node_ids and node_targets, we can even extract actor_nodes and movie nodes....
for i, x in enumerate(node_ids):
    if node_targets[i] == 'actor':
        actor_nodes.append(x)
for i, x in enumerate(node_ids):
    if node_targets[i] == 'movie':
        movie_nodes.append(x)
for i, x in enumerate(node_embeddings):
    if node_targets[i] == 'actor':
        actor_embeddings.append(x)
for i, x in enumerate(node_embeddings):
    if node_targets[i] == 'movie':
        movie_embeddings.append(x)

return actor_nodes, movie_nodes, actor_embeddings, movie_embeddings

```

#lets split the Data

```
actor_nodes, movie_nodes, actor_embeddings, movie_embeddings = data_split(node_ids, node_targets,
```

```
dict: label_map
```

```
print(len(actor_nodes))
```

```
(1 item) {'movie': 0}
```

```
3411
```

## Grader function - 1

```

def grader_actors(data):
    assert(len(data)==3411)
    return True
grader_actors(actor_nodes)

```

```
True
```

## Grader function - 2

```

def grader_movies(data):
    assert(len(data)==1292)
    return True
grader_movies(movie_nodes)

```

```
True
```

```
actor_targets=[ x for x in node_targets if x=='actor']
```



```
movie_targets=[ x for x in node_targets if x=='movie']
```

## Calculating cost1

Cost1 =

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

where N= number of clusters

```
def cost1(graph,number_of_clusters):
```

```
    #calculating cost1'''
```

```
    num= max([len(x) for x in list(nx.connected_components(graph))])
```

```
    Total_Nodes=graph.number_of_nodes()
```

```
    return (1/number_of_clusters)*num/Total_Nodes
```

```
#getting the graph of whole networ using networkx module.....
```

```
import networkx as nx
```

```
dict: label_map
```

```
from networkx.algorithms import bipartite
```

```
(1 item) {'movie': 0}
```

```
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'),('a11','m6'),('a5','m5'),('a10','m8')])
```

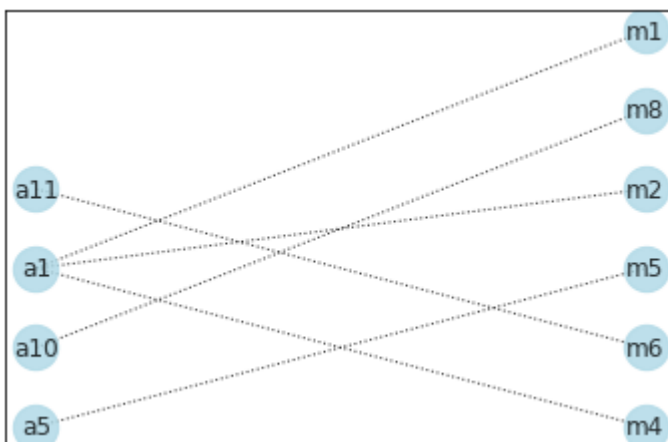
```
l={'a1','a5','a10','a11'};r={'m1','m2','m4','m6','m5','m8'}
```

```
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',alpha=0.8,style='dotted')
```



## Grader function - 3

```

graded_cost1=cost1(graded_graph,3)
def grader_cost1(data):
    assert(data==((1/3)*(4/10)))
    return True
grader_cost1(graded_cost1)

True

```

## 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

```

def cost2(graph,number_of_clusters):

    #calculating cost2'''
    degree = graph.degree()
    nodes = list(graph.nodes())
    unique_nodes = []
    for i in nodes:
        if i not in unique_nodes:
            unique_nodes.append(i)
    summation = 0
    for i in degree:
        if 'a' in i[0]:
            summation+=i[1]
    movie_nodes=0
    for i in unique_nodes:
        if 'm' in i:
            movie_nodes+=1
    return (1/number_of_clusters)*summation/movie_nodes

```

dict: label\_map  
(1 item) {'movie': 0}

## 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)

True

```

## Grouping similar actors

```

number_of_clusters = [3, 5, 10, 30, 50, 100, 200, 500]
cost = []
for cl in number_of_clusters:
    kmeans = KMeans(n_clusters=cl)
    kmeans.fit(actor_embeddings)
    cluster_number_for_data_point = kmeans.labels_
    list_of_all_cluster=[]
    unique = np.unique(cluster_number_for_data_point)
    dict_of_actor_nodes = dict(zip(actor_nodes, cluster_number_for_data_point))
    for number in unique:
        cluster=[]
    for node, cluster_number in dict_of_actor_nodes.items():
        if cluster_number == number:
            cluster.append(node)
        list_of_all_cluster.append(cluster)
cost_1=0
cost_2=0
for cluster_ in list_of_all_cluster:
    G= nx.Graph()
    for actor_node in cluster_:
        sub_graph = nx.ego_graph(B,actor_node)
        G.add_nodes_from(sub_graph.nodes())
        G.add_edges_from(sub_graph.edges())
    cost_1+=cost1(G,cl)
    cost_2+=cost2(G,cl)
print(cost_1*cost_2)
cost.append(cost_1*cost_2)

3764134.4174815724
515083.01906030875
267888.5080381733
25409.597586203927
8000.844928760912
1163.4921000001216
290.8730250000304
23.9970245624989

```

```

dict: label_map
(1 item) {'movie': 0}

```

```
best_cluster=number_of_clusters[cost.index(max(cost))]
```

```

#number of Kmeans Clusters....
algo=KMeans(n_clusters=best_cluster)
algo.fit(actor_embeddings)

```

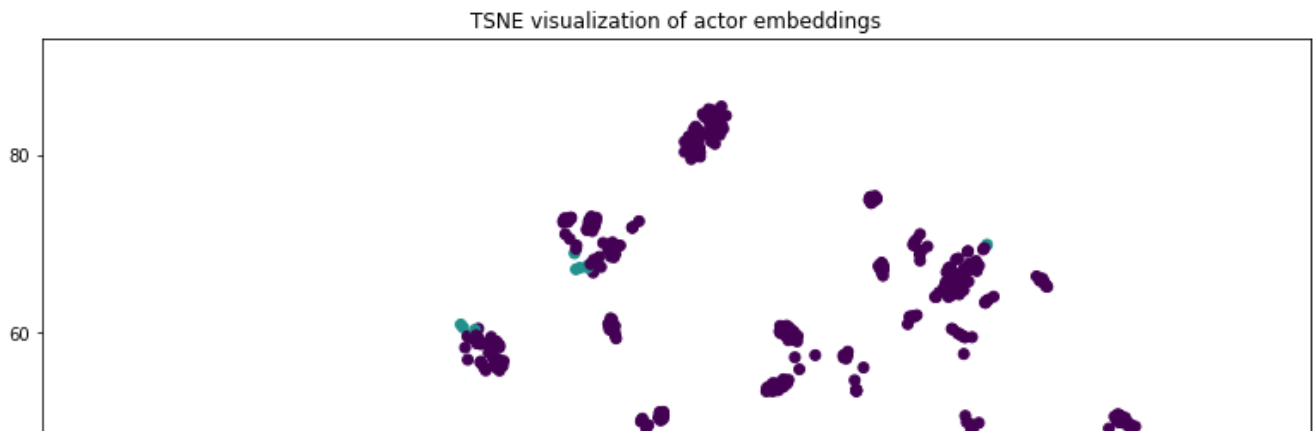
```
KMeans(n_clusters=3)
```

## Displaying similar actor clusters

```
#lets use TSNE with the help of Sklearn
from sklearn.manifold import TSNE
transform = TSNE
trans = transform(n_components=2)
actor_embeddings_2d = trans.fit_transform(actor_embeddings)
label_map = { l: i for i, l in enumerate(np.unique(actor_targets))}
actor_colours = [ label_map[target] for target in actor_targets]
plt.figure(figsize=(20,16))
plt.axes().set(aspect="equal")
plt.scatter(actor_embeddings_2d[:,0],actor_embeddings_2d[:,1],c=algo.predict(actor_embeddings_2d))
plt.title('{} visualization of actor embeddings'.format(transform.__name__))
plt.show()
```



```
dict: label_map
(1 item) {'movie': 0}
```



### Grouping similar movies

```
cluster_list=[3,5,10,30,50,100,200,500]
Cost_movies=[]
for cluster in cluster_list:
    algo_m=KMeans(n_clusters=cluster)
    algo_m.fit(movie_embeddings)
    label_m=algo_m.labels_
    dic=dict(zip(movie_nodes,label_m))
    c1=0
    c2=0
    for i in label_m:
        dict: label_map
        ac_node = [k for k,v in dic.items() if v == i] (1 item) {'movie': 0}
        G1=nx.Graph()
        for n in range(len(ac_node)):
            sub_graph1 = nx.ego_graph(A,node_ids[n])
            G1.add_nodes_from(sub_graph1.nodes)
            G1.add_edges_from(sub_graph1.edges())
            c1+=cost1(G1,cluster)
            c2+=cost2(G1,cluster)
    print(c1*c2)
    Cost_movies.append(c1*c2)

1164804.88819142
357792.63769194606
77616.87231410101
6599.677377353273
2036.2201280541808
451.71919964291226
97.10512671425177
10.927438307722747
```

```
best_cluster=cluster_list[Cost_movies.index(max(Cost_movies))]
```

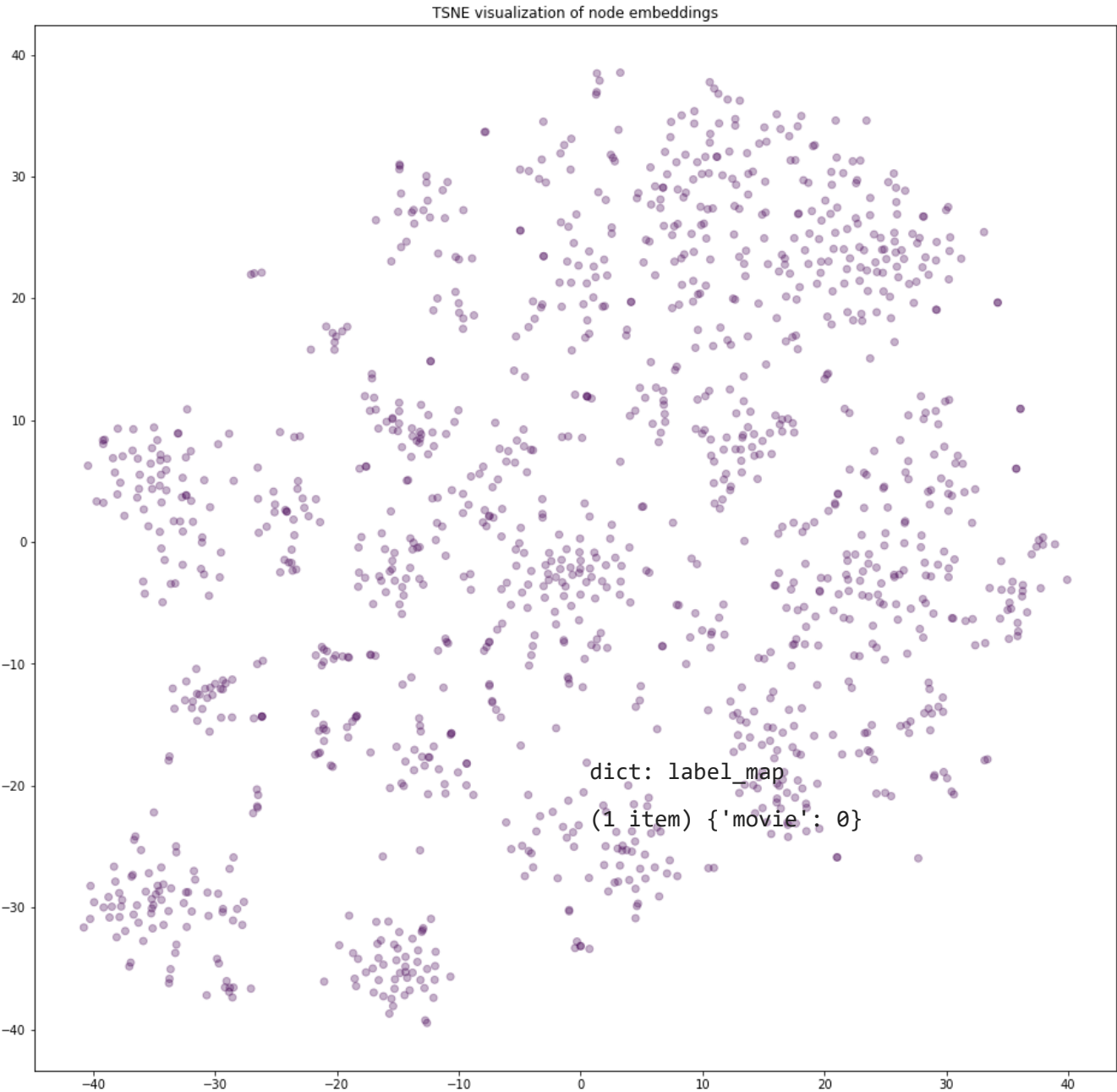
```
kmeans=KMeans(n_clusters=best_cluster)
kmeans.fit(movie_embeddings)
```

```
KMeans(n_clusters=3)
```

## Displaying similar movie clusters

```
from sklearn.manifold import TSNE
transform = TSNE
trans_ = transform(n_components=2)
movie_embeddings_2d = trans_.fit_transform(movie_embeddings)
import numpy as np
# drawing the points
label_map = { l: i for i, l in enumerate(np.unique(movie_targets))}
node_colours = [ label_map[target] for target in movie_targets]
plt.figure(figsize=(20,16))
plt.axes().set(aspect="equal")
plt.scatter(movie_embeddings_2d[:,0],movie_embeddings_2d[:,1],c=node_colours, alpha=0.3)
plt.title('{} visualization of node embeddings'.format(transform.__name__))
plt.show()
```

```
dict: label_map
(1 item) {'movie': 0}
```



✓ 1s completed at 12:06 PM

