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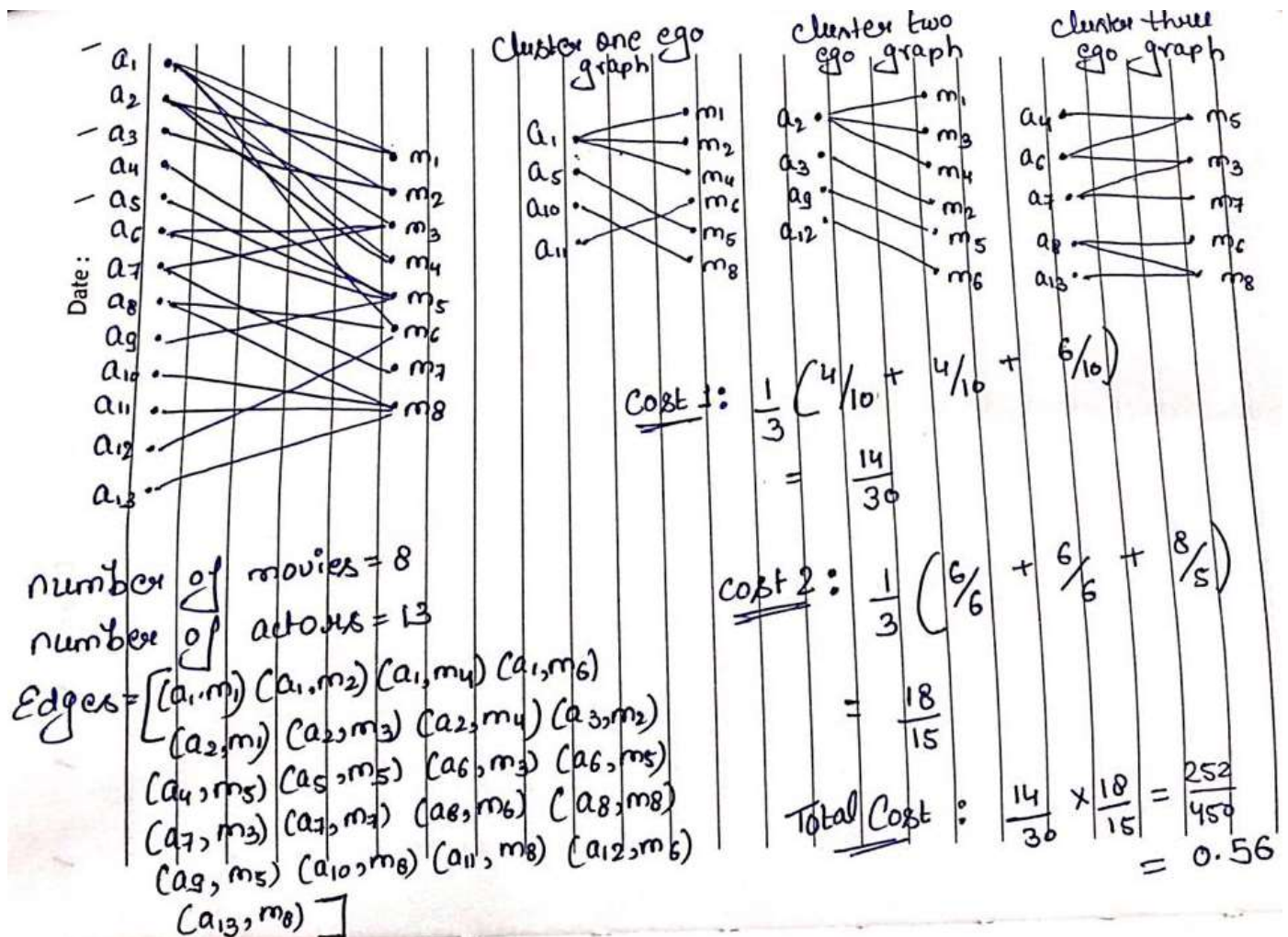
**Every Grader function has to return True.**

**Please check [clustering assignment helper functions (https://drive.google.com/file/d/1V29KhKo3YnckMX32treEgdtH5r90DIjU/view?usp=sharing)](https://drive.google.com/file/d/1V29KhKo3YnckMX32treEgdtH5r90DIjU/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)](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 =
   $\frac{1}{N} \sum_{\text{each cluster i}} \frac{(\text{number of nodes in the largest connected component in the graph with the actor nodes and its movie neighbours in clus}}{(\text{total number of nodes in that cluster i})}$
   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

Cluster one ego graph

Cluster two ego graph

Cluster three ego graph

number of movies = 8

number of actors = 13

$$Edges = \begin{bmatrix} (a_1,m_1) \; (a_1,m_2) \; (a_1,m_4) \; (a_1,m_6) \\ (a_2,m_1) \; (a_2,m_3) \; (a_2,m_4) \; (a_3,m_2) \\ (a_4,m_5) \; (a_5,m_5) \; (a_6,m_3) \; (a_6,m_5) \\ (a_7,m_3) \; (a_7,m_7) \; (a_8,m_6) \; (a_8,m_8) \\ (a_9,m_5) \; (a_{10},m_8) \; (a_{11},m_8) \; (a_{12},m_6) \\ (a_{13},m_8) \end{bmatrix}$$

Cost 1: $\frac{1}{3}\left(\frac{4}{10} + \frac{4}{10} + \frac{6}{10}\right)$

$= \frac{14}{30}$

Cost 2: $\frac{1}{3}\left(\frac{6}{6} + \frac{6}{6} + \frac{8}{5}\right)$

$= \frac{18}{15}$

Total Cost: $\frac{14}{30} \times \frac{18}{15} = \frac{252}{450}$

$= 0.56$

# 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 =
$$\frac{1}{N} \sum_{\text{each cluster i}} \frac{\text{(number of nodes in the largest connected component in the graph with the movie nodes and its actor neighbours in clus}}{\text{(total number of nodes in that cluster i)}}$$
where N= number of clusters
(Write your code in def cost1())

3. Cost2 = $\frac{1}{N} \sum_{\text{each cluster i}} \frac{\text{(sum of degress of movie nodes in the graph with the movie nodes and its actor neighbours in cluster i)}}{\text{(number of unique actor nodes in the graph with the movie nodes and its actor neighbours in cluster i)}}$
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_algorith(clusters=number_of_clusters)
    # you will be passing a matrix of size N*d where N number of actor nodes a
```

```
nd 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
```

In [1]:

```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
# you need to have tensorflow
from stellargraph.data import UniformRandomMetaPathWalk
from stellargraph import StellarGraph
```

In [2]:

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

In [3]:

```python
edges = [tuple(x) for x in data.values.tolist()]
```
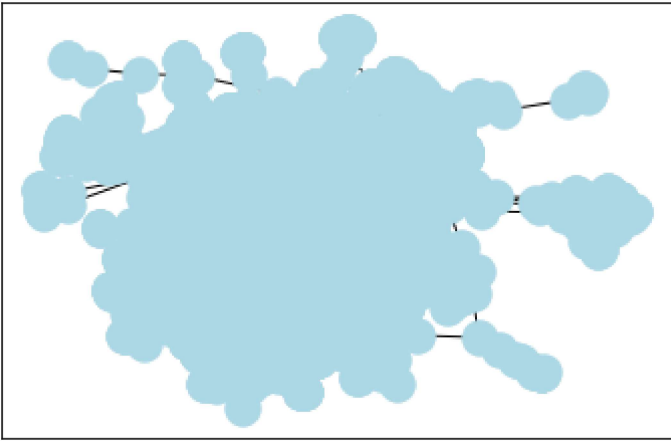
In [4]:

```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 [5]:

```python
nx.draw_networkx(B,node_color='lightblue',with_labels=False)
```



In [5]:

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

In [6]:

```python
print("number of nodes", A.number_of_nodes())
print("number of edges", A.number_of_edges())
```
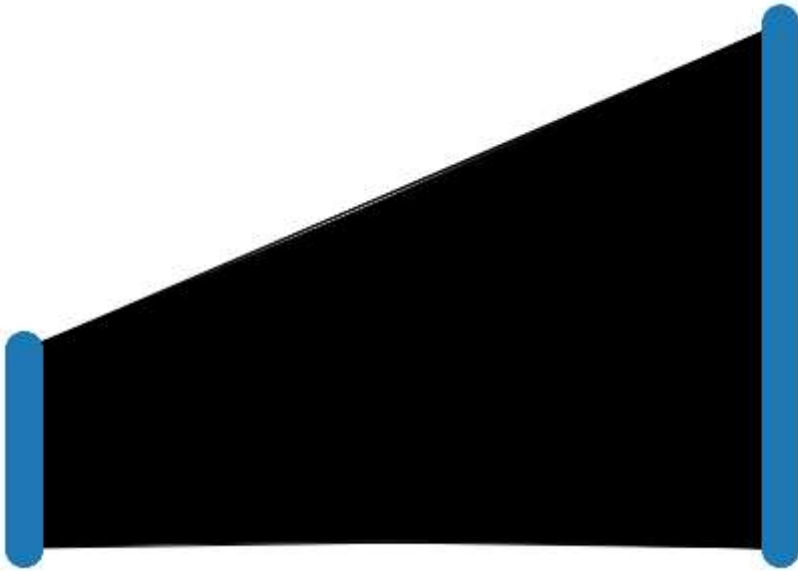
```
number of nodes 4703
number of edges 9650
```

In [8]:

```
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=False)
plt.show()
```

In [7]:

```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 [8]:

```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 [9]:

```python
from gensim.models import Word2Vec
model = Word2Vec(walks, vector_size=128, window=5)
```

In [10]:

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

Out[10]:

```
(4703, 128)
```

In [11]:

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

In [12]:

```
print(node_embeddings.shape)
```

```
(4703, 128)
```

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

In [13]:

```
a_split(node_ids,node_targets,node_embeddings):
    In this function, we will split the node embeddings into actor_embeddings , movie_embeddings
    or_nodes,movie_nodes=[],[]
    or_embeddings,movie_embeddings=[],[]
    or_targets,movie_targets=[],[]
    plit the node_embeddings into actor_embeddings,movie_embeddings based on node_ids
    y using node_embedding and node_targets, we can extract actor_embedding and movie embedding
    y using node_ids and node_targets, we can extract actor_nodes and movie nodes
    i,item in enumerate(node_ids):
    if 'm' in item:
        movie_nodes.append(item)
        movie_embeddings.append(node_embeddings.tolist()[i])
        movie_targets.append(node_targets[i])
    else:
        actor_nodes.append(item)
        actor_embeddings.append(node_embeddings.tolist()[i])
        actor_targets.append(node_targets[i])
    urn actor_nodes,movie_nodes,actor_embeddings,movie_embeddings,actor_targets,movie_targets

    odes,movie_nodes,actor_embeddings,movie_embeddings,actor_targets,movie_targets = data_split(
```

Grader function - 1

In [14]:

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

Out[14]:

```
True
```

Grader function - 2

In [15]:

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

Out[15]:

True

## Calculating cost1

Cost1 =

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

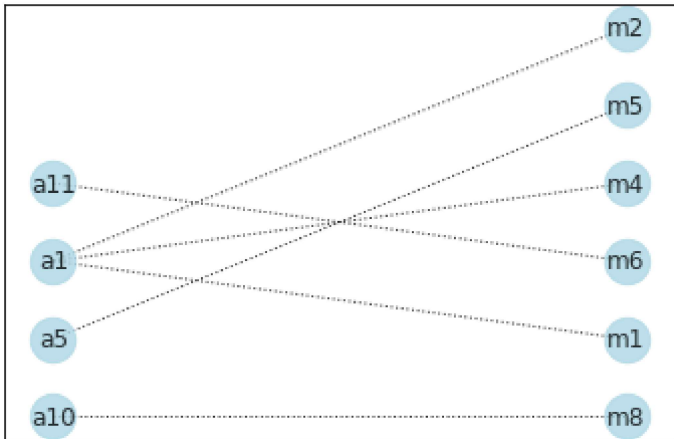where N= number of clusters

In [16]:

```python
def find_clusters(number_of_clusters,lblst,nodes):
    cluster_lst = []
    for i in range(number_of_clusters):
        cluster = []
        for j,val in enumerate(lblst):
            if val == i:
                cluster.append(nodes[j])
        cluster_lst.append(cluster)
    return cluster_lst
```

In [17]:

```python
def cost1(graph,number_of_clusters):
    '''In this function, we will calculate cost1'''
    lrgst_conn_comp_nodes = len(max(nx.connected_components(graph), key=len))
    Clust_cost = lrgst_conn_comp_nodes/len(graph.nodes())
    return Clust_cost/number_of_clusters
```

In [18]:

```python
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 node attribute
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'),(
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,s
```



Grader function - 3

In [19]:

```python
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[19]:

True

## Calculating cost2

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

In [20]:

```python
def cost2(graph,number_of_clusters):
    '''In this function, we will calculate cost1'''
    total_nodes = graph.nodes()
    uniq_mve_nodes = 0
    for node in total_nodes:
        if 'm' in node:
            uniq_mve_nodes+=1
    Clust_cost = sum(dict(graph.degree()).values())/(2*uniq_mve_nodes)
    return Clust_cost/number_of_clusters
```

In [21]:

```python
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[21]:

True

In [22]:

```python
from sklearn.cluster import KMeans
cost_dict = {}
for number_of_clusters in [3, 5, 10, 30, 50, 100, 200, 500]:
    algo = KMeans(n_clusters=number_of_clusters)
    algo.fit(actor_embeddings)
    lblst = algo.labels_.tolist()
    actor_clusters = find_clusters(number_of_clusters,lblst,actor_nodes)
    cost1_val_lst = []
    cost2_val_lst = []
    for cluster in actor_clusters:
        subgrph_lst = []
        for node in cluster:
            sub_graph1=nx.ego_graph(B,node)
            subgrph_lst.append(sub_graph1)
        comb_graph = nx.compose_all(subgrph_lst)
        cost1_val = cost1(comb_graph,number_of_clusters)
        cost2_val = cost2(comb_graph,number_of_clusters)
        cost1_val_lst.append(cost1_val)
        cost2_val_lst.append(cost2_val)
    metric_Cost = sum(cost1_val_lst)*sum(cost2_val_lst)
    cost_dict[number_of_clusters] = metric_Cost
max_cost_num_clust = sorted(cost_dict.items(),key=lambda x: x[1],reverse=True)[0][0]
print(max_cost_num_clust)
```

3

Grader function - 4

Grouping similar actors

In [23]:

```python
algo = KMeans(n_clusters=3)
algo.fit(actor_embeddings)
lblst = algo.labels_.tolist()
actor_clusters = find_clusters(number_of_clusters,lblst,actor_nodes)
```

Displaying similar actor clusters

In [29]:

```python
import numpy as np
from sklearn.manifold import TSNE

transform = TSNE #PCA

trans = transform(n_components=2)
actor_embeddings_2d = trans.fit_transform(actor_embeddings)

# draw the points

#node_colours = np.array(["red","green","blue"]) # best #of clusters are 3. So, 3 colors ar
label_map = { l: i for i, l in enumerate(np.unique(lblst))}
node_colours = np.array([ label_map[target] for target in lblst])

plt.figure(figsize=(20,16))
plt.axes().set(aspect="equal")
plt.scatter(actor_embeddings_2d[:,0],
            actor_embeddings_2d[:,1],
            c=node_colours, alpha=0.3)
plt.title('{} visualization of actor embeddings'.format(transform.__name__))

plt.show()
```

TSNE visualization of actor embeddings

In [30]:

```python
cost_dict = {}
for number_of_clusters in [3, 5, 10, 30, 50, 100, 200, 500]:
    algo = KMeans(n_clusters=number_of_clusters)
    algo.fit(movie_embeddings)
    lblst = algo.labels_.tolist()
    movie_clusters = find_clusters(number_of_clusters,lblst,movie_nodes)
    cost1_val_lst = []
    cost2_val_lst = []
    for cluster in movie_clusters:
        subgrph_lst = []
        for node in cluster:
            sub_graph1=nx.ego_graph(B,node)
            subgrph_lst.append(sub_graph1)
        comb_graph = nx.compose_all(subgrph_lst)
        cost1_val = cost1(comb_graph,number_of_clusters)
        cost2_val = cost2(comb_graph,number_of_clusters)
        cost1_val_lst.append(cost1_val)
        cost2_val_lst.append(cost2_val)
    metric_Cost = sum(cost1_val_lst)*sum(cost2_val_lst)
    cost_dict[number_of_clusters] = metric_Cost
max_cost_num_clust = sorted(cost_dict.items(),key=lambda x: x[1],reverse=True)[0][0]
print(max_cost_num_clust)
```

30

## Grouping similar movies

In [31]:

```python
algo = KMeans(n_clusters=50)
algo.fit(movie_embeddings)
lblst = algo.labels_.tolist()
movie_clusters = find_clusters(number_of_clusters,lblst,movie_nodes)
print(movie_clusters[0])
```

```
['m1357', 'm69', 'm1366', 'm941', 'm1367', 'm833', 'm825', 'm669', 'm1361',
 'm1380', 'm65', 'm831', 'm832', 'm1228', 'm837', 'm1319', 'm1375', 'm1369',
 'm1358', 'm1323', 'm1379', 'm1350']
```
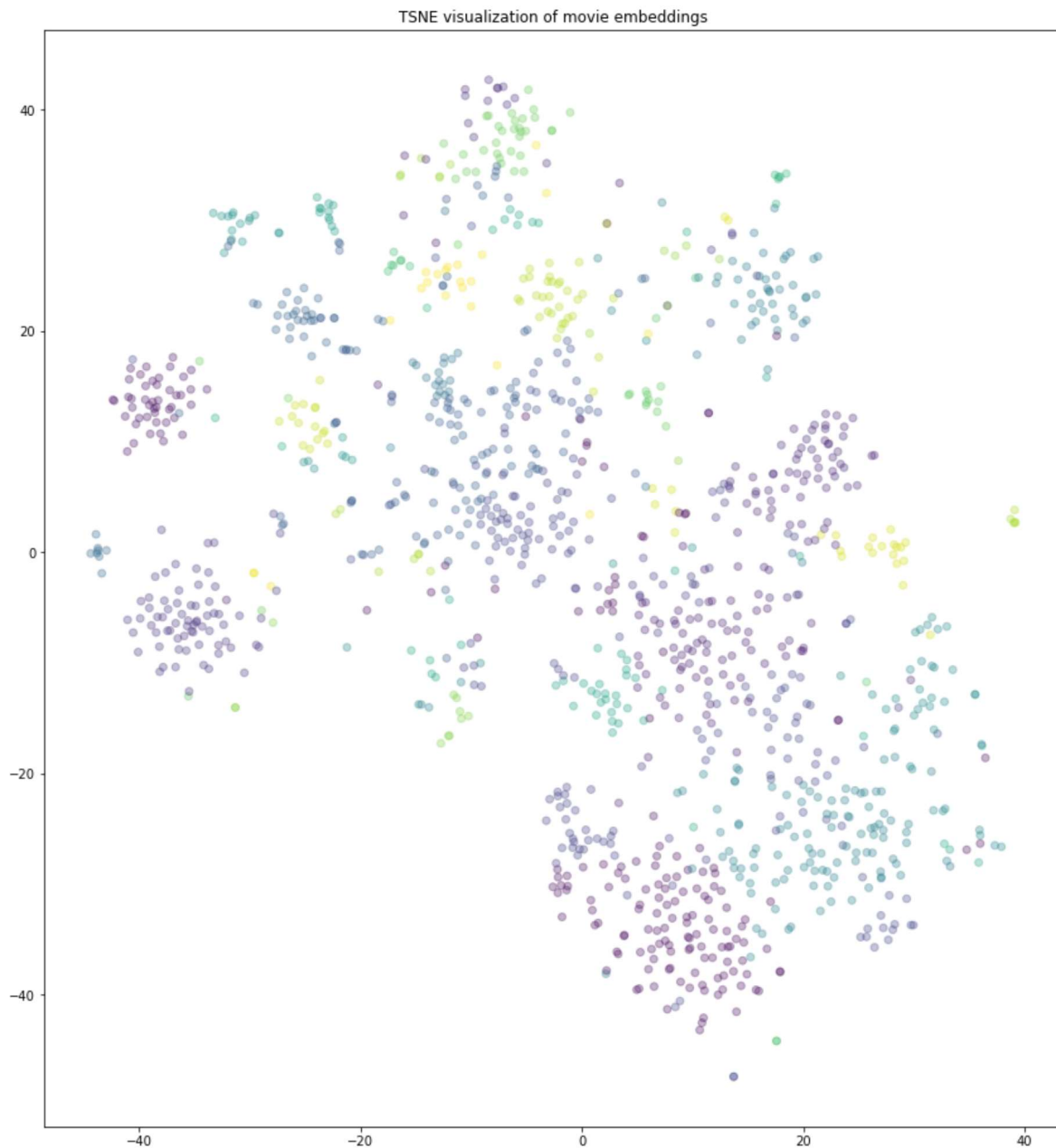
## Displaying similar movie clusters

In [33]:

```python
movie_embeddings_2d = trans.fit_transform(movie_embeddings)

# draw the points
label_map = { l: i for i, l in enumerate(np.unique(lblst))}
node_colours = np.array([ label_map[target] for target in lblst])

plt.figure(figsize=(20,16))
plt.axes().set(aspect="equal")
plt.scatter(movie_embeddings_2d[:,0].flatten(),
            movie_embeddings_2d[:,1].flatten(),
            c=node_colours, alpha=0.3)
plt.title('{} visualization of movie embeddings'.format(transform.__name__))

plt.show()
```



TSNE visualization of movie embeddings

In [ ]: