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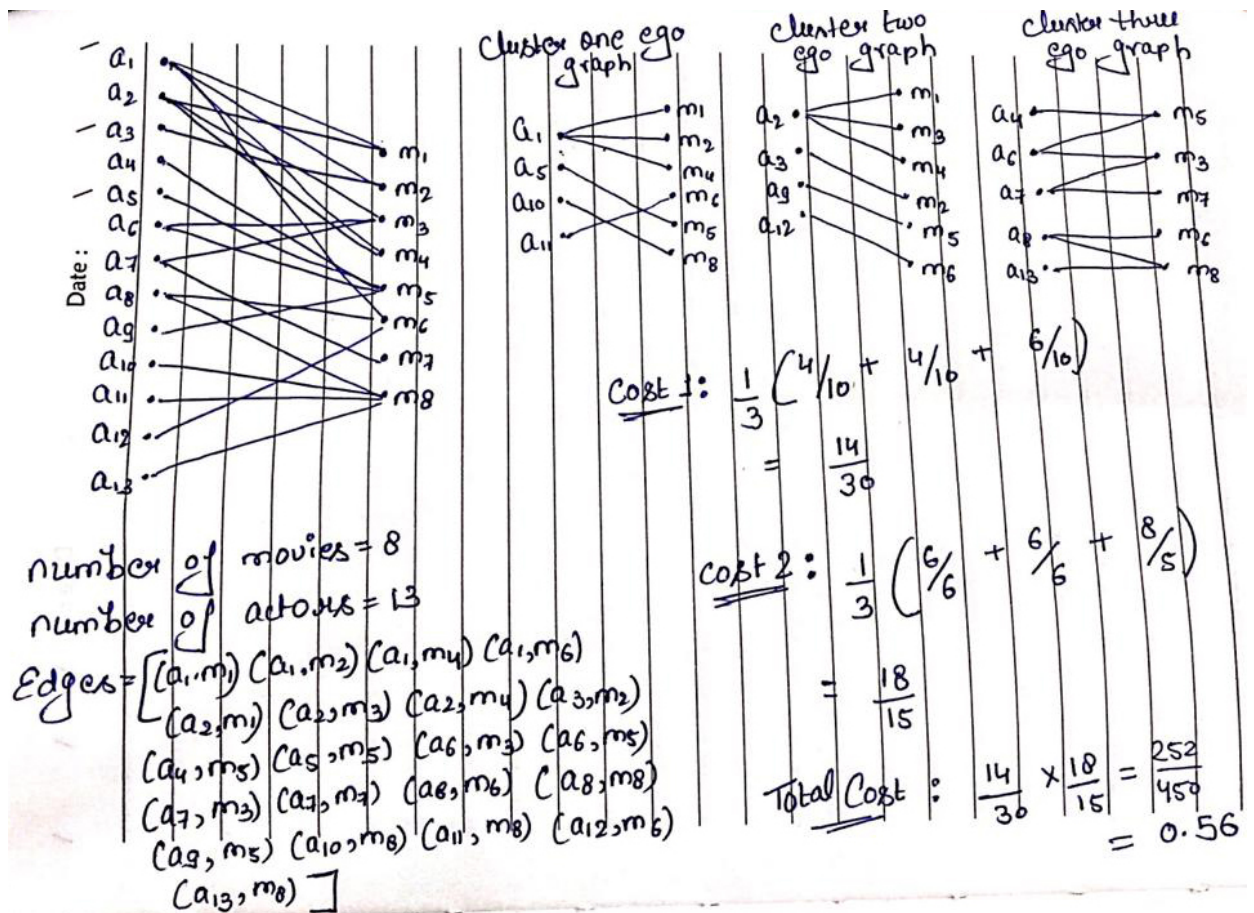
**Every Grader function has to return True.**

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

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

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 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
```

In [2]:
```python
!pip install networkx==2.3
```

```
Requirement already satisfied: networkx==2.3 in d:\anaconda3\lib\site-packages (2.3)
Requirement already satisfied: decorator>=4.3.0 in d:\anaconda3\lib\site-packages (from
networkx==2.3) (4.4.2)
```

In [3]:
```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 [4]:
```python
data=pd.read_csv('movie_actor_network.csv', index_col=False, names=['movie','actor'])
```

In [5]:
```python
edges = [tuple(x) for x in data.values.tolist()]
```

In [6]:
```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 [7]:
```python
A = list(nx.connected_component_subgraphs(B))[0]
```

In [8]:
```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 [9]:
```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 [10]:
```
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 [11]:
```
# 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 [12]:
```
from gensim.models import Word2Vec
model = Word2Vec(walks, size=128, window=5)
```

In [13]:
```
model.wv.vectors.shape  # 128-dimensional vector for each node in the graph
```

Out[13]: (4703, 128)

In [14]:
```
# 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 embed
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', 'movie', 'actor', 'movie', 'actor', 'movie']
```

In [15]:
```
def data_split(node_ids,node_targets,node_embeddings):
    '''In this function, we will split the node embeddings into actor_embeddings , movi
    actor_nodes,movie_nodes=[],[]
    actor_embeddings,movie_embeddings=[],[]
    # split the node_embeddings into actor_embeddings,movie_embeddings based on node_id
```

```
        # By using node_embedding and node_targets, we can extract actor_embedding and movi
        # By using node_ids and node_targets, we can extract actor_nodes and movie nodes

        for i,j in enumerate(node_targets):
            if j=="actor":
                actor_nodes.append(node_ids[i])
            else:
                movie_nodes.append(node_ids[i])

        for i,j in enumerate(node_targets):
            if j=="actor":
                actor_embeddings.append(node_embeddings[i])
            else:
                movie_embeddings.append(node_embeddings[i])

        return actor_nodes,movie_nodes,actor_embeddings,movie_embeddings

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

### Grader function - 1

In [16]:
```python
def grader_actors(data):
    assert(len(data)==3411)
    return True
grader_actors(actor_nodes)
```

Out[16]:  True

### Grader function - 2

In [17]:
```python
def grader_movies(data):
    assert(len(data)==1292)
    return True
grader_movies(movie_nodes)
```

Out[17]:  True

### Calculating cost1

Cost1 =

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

where N= number of clusters

In [18]:
```python
def cost1(graph,number_of_clusters):
    '''In this function, we will calculate cost1'''

    #https://networkx.org/documentation/stable//reference/algorithms/generated/networkx
    Num_of_nodes_in_largest_cc = max(nx.connected_component_subgraphs(graph), key=len).

    summation = Num_of_nodes_in_largest_cc/graph.number_of_nodes()

    cost1= summation/number_of_clusters
    return cost1
```

```python
import networkx as nx
```

```
In [19]:   from networkx.algorithms import bipartite
           graded_graph= nx.Graph()
           graded_graph.add_nodes_from(['a1','a5','a10','a11'], bipartite=0) # Add the node attrib
           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
```

### Grader function - 3

```
In [20]:   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[20]:   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 [21]:   def cost2(graph,number_of_clusters):
               '''In this function, we will calculate cost1'''

               sum_of_a_deg = 0
               for i,j in enumerate(graph.degree()):
                   if 'a' in j[0]:
                       sum_of_a_deg += j[1]

               total_m = 0
               for nod in graph.nodes():
                   if 'm' in nod:
                       total_m+=1

               cost2= (sum_of_a_deg/total_m)/number_of_clusters

               return cost2
```

### Grader function - 4

```
In [22]:   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[22]:   True

### Grouping similar actors

```
In [64]:   # 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 and
```

```
#            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
#                 cost2=cost2(graph1)+cost2(graph2)+cost2(graph3)
#            computer the metric Cost = Cost1*Cost2
#       return number_of_clusters which have maximum Cost

from sklearn.cluster import KMeans

costs = []
No_cluster = [3, 5, 10, 30, 50, 100, 200, 500]
for number_of_clusters in No_cluster:
    print("Number of Clusters : " , number_of_clusters)
    kmeans = KMeans(n_clusters=number_of_clusters, random_state=0)
    kmeans.fit(actor_embeddings)
    lbls = kmeans.labels_

    mapping_dic = {}
    for i in range(len(lbls)):
        mapping_dic[actor_nodes[i]] = lbls[i]

    c1 = 0
    c2 = 0
    for i in range(number_of_clusters):
        G1=nx.Graph()
        for k,v in mapping_dic.items():
            if(v==i):
                sub_graph = nx.ego_graph(B,k)
                G1.add_nodes_from(sub_graph.nodes)
                G1.add_edges_from(sub_graph.edges())
        c1 += cost1(G1,number_of_clusters)
        c2 += cost2(G1,number_of_clusters)
    costs.append(c1*c2)

idx_maxCost = costs.index(max(costs))

print ("Number of clusters that gives the highest score : {} , Score : {} ".format(No_c
```

```
Number of Clusters :  3
Number of Clusters :  5
Number of Clusters :  10
Number of Clusters :  30
Number of Clusters :  50
Number of Clusters :  100
Number of Clusters :  200
Number of Clusters :  500
Number of clusters that gives the highest score : 3 , Score : 3.8087351005437453
```

Displaying similar actor clusters

```
In [68]:  kmeans = KMeans(n_clusters=3, random_state=0)
          kmeans.fit(actor_embeddings)
          lbls = kmeans.labels_

          from sklearn.manifold import TSNE
          transform = TSNE #PCA

          trans = transform(n_components=2)
          node_embeddings_2d = trans.fit_transform(actor_embeddings)
```
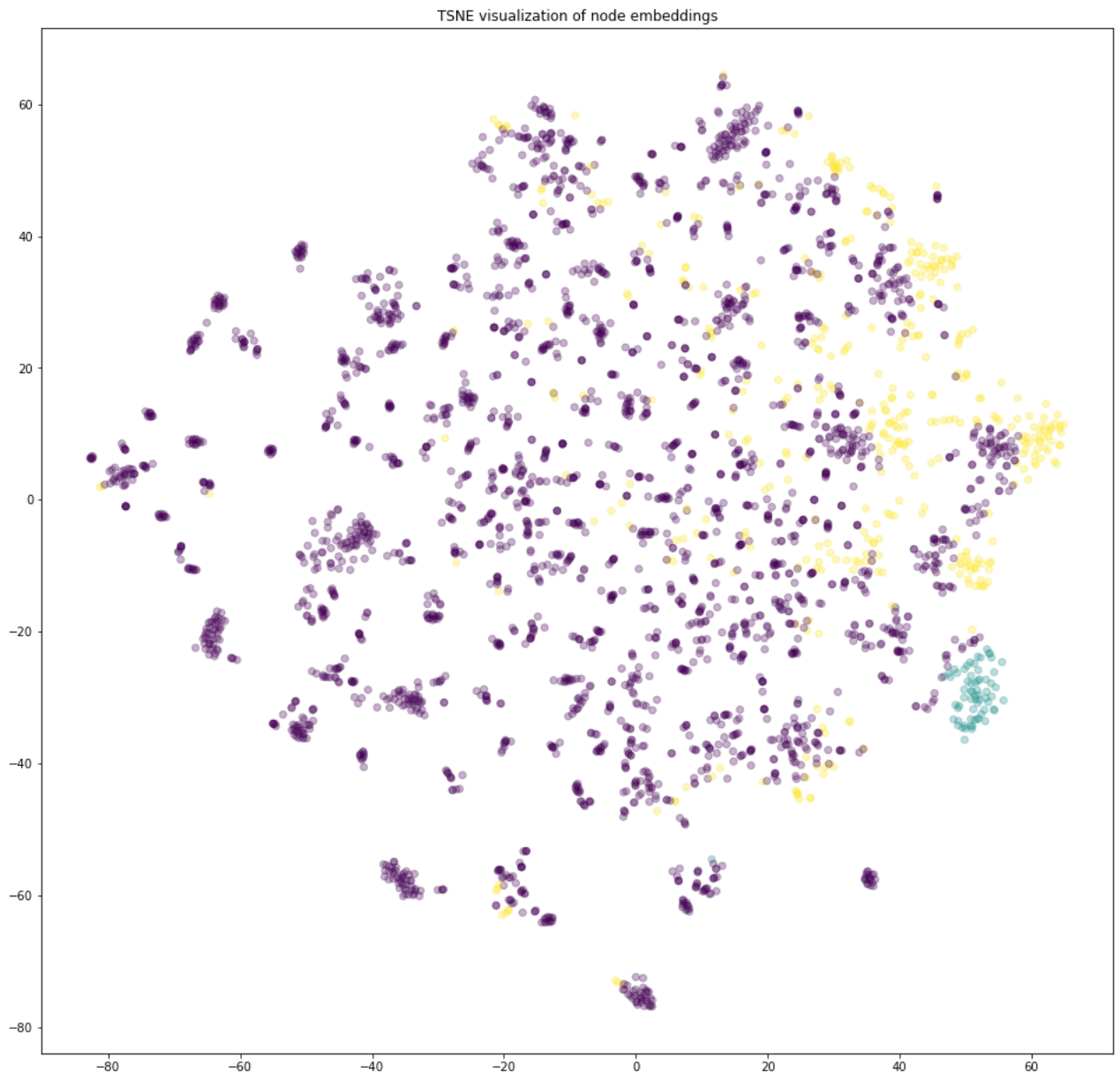
In [76]:
```python
mapping_dic = {}
for i in range(len(lbls)):
    mapping_dic[actor_nodes[i]] = lbls[i]
```

In [79]:
```python
import numpy as np
# draw the points

label_map = mapping_dic
node_colours = [ label_map[target] for target in actor_nodes]

plt.figure(figsize=(20,16))
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()
```



TSNE visualization of node embeddings

Grouping similar movies

In [80]:
```python
def cost1(graph,number_of_clusters):
    '''In this function, we will calculate cost1'''

    #https://networkx.org/documentation/stable//reference/algorithms/generated/networkx
    Num_of_nodes_in_largest_cc = max(nx.connected_component_subgraphs(graph), key=len).

    summation = Num_of_nodes_in_largest_cc/graph.number_of_nodes()

    cost1= summation/number_of_clusters
    return cost1


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

    sum_of_a_deg = 0

    for i,j in enumerate(graph.degree()):

        if 'm' in j[0]:
            sum_of_a_deg += j[1]

    total_m = 0
    for nod in graph.nodes():
        if 'a' in nod:
            total_m+=1

    cost2= (sum_of_a_deg/total_m)/number_of_clusters

    return cost2

from sklearn.cluster import KMeans

costs = []
No_cluster = [3, 5, 10, 30, 50, 100, 200, 500]
for number_of_clusters in No_cluster:
    print("Number of Clusters : " , number_of_clusters)
    kmeans = KMeans(n_clusters=number_of_clusters, random_state=0)
    kmeans.fit(movie_embeddings)
    lbls = kmeans.labels_

    mapping_dic = {}
    for i in range(len(lbls)):
        mapping_dic[movie_nodes[i]] = lbls[i]

    c1 = 0
    c2 = 0
    for i in range(number_of_clusters):
        G1=nx.Graph()
        for k,v in mapping_dic.items():
            if(v==i):
                sub_graph = nx.ego_graph(B,k)
                G1.add_nodes_from(sub_graph.nodes)
                G1.add_edges_from(sub_graph.edges())
        c1 += cost1(G1,number_of_clusters)
        c2 += cost2(G1,number_of_clusters)
    costs.append(c1*c2)

idx_maxCost = costs.index(max(costs))

print ("Number of clusters that gives the highest score : {} , Score : {} ".format(No_c
```

```
Number of Clusters :  3
Number of Clusters :  5
Number of Clusters :  10
Number of Clusters :  30
Number of Clusters :  50
Number of Clusters :  100
Number of Clusters :  200
Number of Clusters :  500
Number of clusters that gives the highest score : 3 , Score : 2.6789381237318244
```

## Displaying similar movie clusters

In [83]:
```python
kmeans = KMeans(n_clusters=3, random_state=0)
kmeans.fit(movie_embeddings)
lbls = kmeans.labels_

from sklearn.manifold import TSNE
transform = TSNE #PCA

trans = transform(n_components=2)
node_embeddings_2d = trans.fit_transform(movie_embeddings)

mapping_dic = {}
for i in range(len(lbls)):
    mapping_dic[movie_nodes[i]] = lbls[i]

import numpy as np
# draw the points

label_map = mapping_dic
node_colours = [ label_map[target] for target in movie_nodes]

plt.figure(figsize=(20,16))
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()
```

TSNE visualization of node embeddings