# AML_HW4_Write up

Students

- Scarlett Huang (sh2557) CM
- Zihan Zhang (zz698) ORIE

Submission: Late

## Programming Exercises

## Question 1

### 1.(a)

```python
# import modules

import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
%matplotlib inline
import math
```

```python
# load .npy file

dw_matrix = np.load('data/science2k-doc-word.npy')

print (dw_matrix)
print (len(dw_matrix))
print (type(dw_matrix))
print (dw_matrix.shape)
```

```
[[-0.2521619 -0.2521619  9.36371   ... -0.2521619 -0.2521619
  -0.2521619]
 [-0.2875293 -0.2875293  8.229864  ... -0.2875293 -0.2875293
  -0.2875293]
 [-0.3634041 -0.3634041  9.252468  ... -0.3634041 -0.3634041
  -0.3634041]
 ...
 [10.04846    -0.7713402  9.132197  ... -0.7713402 -0.7713402
  -0.7713402]
 [11.00702    -0.8423803 10.38288   ... -0.8423803 -0.8423803
  -0.8423803]
 [ 9.710337   -0.7527946  9.556191  ... -0.7527946 -0.7527946
  -0.7527946]]
1373
<class 'numpy.ndarray'>
(1373, 5476)
```

```python
# a-1: implement kmeans

from numpy import *
import matplotlib.pyplot as plt

# Compute Euclidean distance
def euclidean_distance(vec1, vec2):
  eu_distance = sqrt(sum(power(vec2 - vec1, 2)))
  return eu_distance

# init centroids with random samples
def init_centroids(data, k):
  num_samples, dim = data.shape
  centroids = zeros((k, dim))
  for i in range(k):
    idx = int(random.uniform(0, num_samples))
    centroids[i, :] = data[idx, :]
  return centroids

# kmeans
def kmeans(data, k):
  num_samples = data.shape[0]

```

```python
24      doc_cluster = mat(zeros((num_samples, 2)))
25      cls_changed = True
26
27    # init centroids
28    centroids = init_centroids(data, k)
29
30    while cls_changed:
31      cls_changed = False
32      for i in range(num_samples):
33        min_dist  = 100000.0
34        min_idx = 0
35        # for each centroid, find the centroid who is closest
36        for j in range(k):
37          distance = euclidean_distance(centroids[j, :], data[i,
    :])
38          if distance < min_dist:
39            min_dist  = distance
40            min_idx = j
41
42        # update its cluster
43        if doc_cluster[i, 0] != min_idx:
44          cls_changed = True
45          doc_cluster[i, :] = min_idx, min_dist**2
46
47      # update centroids
48      for j in range(k):
49        points_cls = data[nonzero(doc_cluster[:, 0].A == j)[0]]
50        centroids[j, :] = mean(points_cls, axis = 0)
51
52    return centroids, doc_cluster
53
54  # plot the cluster
55  def show_cluster(data, k, centroids, doc_cluster):
56    num_samples, dim = data.shape
57    if dim != 5476:
58      return 1
59
60    marks = ['or', 'ob', 'og', 'ok', '^r', '+r', 'sr', 'dr',
    '<r', 'pr','or', 'ob', 'og', 'ok', '^r', '+r', 'sr', 'dr',
    '<r', 'pr']
61    if k > len(marks):
62      return 1
```
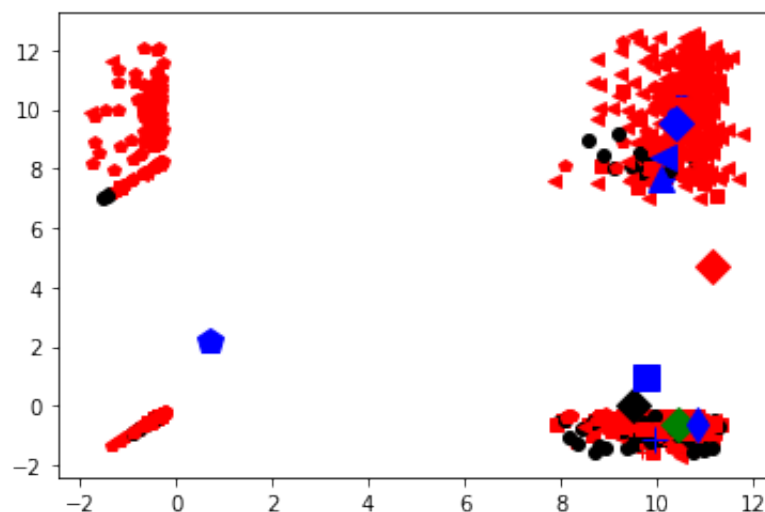
```
63
64    # plot all samples
65    for i in range(num_samples):
66      mark_idx = int(doc_cluster[i, 0])
67      plt.plot(data[i, 0], data[i, 1], marks[mark_idx])
68
69    marks = ['Dr', 'Db', 'Dg', 'Dk', '^b', '+b', 'sb', 'db',
      '<b', 'pb','Dr', 'Db', 'Dg', 'Dk', '^b', '+b', 'sb', 'db',
      '<b', 'pb']
70    # draw the centroids
71    for i in range(k):
72      plt.plot(centroids[i, 0], centroids[i, 1], marks[i],
    markersize = 12)
73
74    plt.show()
```

```
1  # Run the model
2  ## clustering, k=10
3  data = mat(dw_matrix)
4  k = 10
5  centroids, doc_cluster = kmeans(data, k)
6
7  ## plot
8  show_cluster(data, k, centroids, doc_cluster)
```
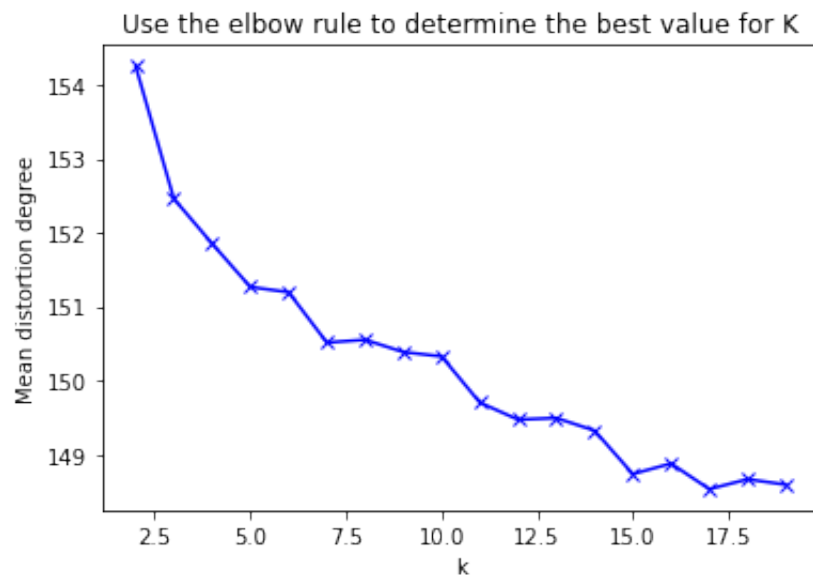


```
1  # a-2: Select the best value for k
2
3  from sklearn.cluster import KMeans
```

```
4   from sklearn import metrics
5   from scipy.spatial.distance import cdist
6   import matplotlib.pyplot as plt
7
8   K = range(2, 20)
9   meandistortions = []
10
11  X = mat(dw_matrix)
12
13  for k in K:
14      kmeans = KMeans(n_clusters=k)
15      kmeans.fit(X)
16      meandistortions.append(sum(np.min(cdist(X,
    kmeans.cluster_centers_, 'euclidean'), axis=1)) / X.shape[0])
17
18  plt.plot(K, meandistortions, 'bx-')
19  plt.xlabel('k')
20  plt.ylabel(u'Mean distortion degree')
21  plt.title(u'Use the elbow rule to determine the best value for
    K');
22  plt.show()
```


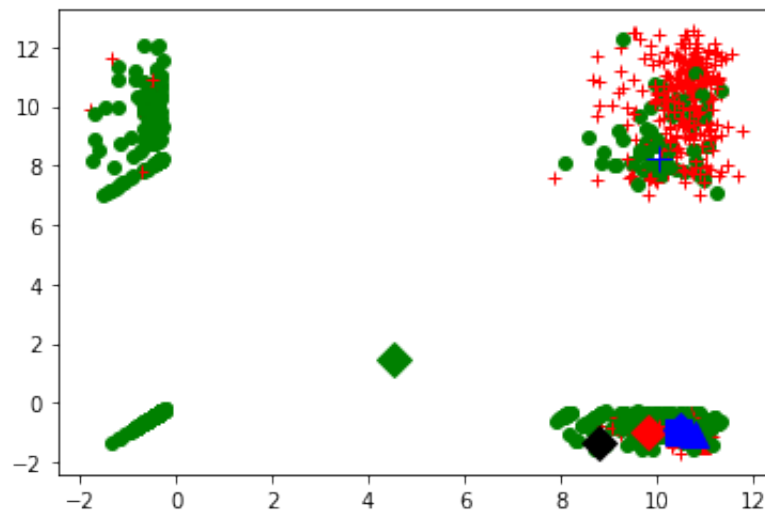Use the elbow rule to determine the best value for K

```
1   # a-3: Run the model
2
3   ## clustering
4   ## set k=7, because according to elbow rule, we can see that in
    the process of K value increasing, the K value corresponding to
    the position where the improvement effect of average distortion
    degree decreases the most is 7.
5
6   data = mat(dw_matrix)
7   k = 7
8   centroids, doc_cluster = kmeans(data, k)
9
10  ## plot
11  show_cluster(data, k, centroids, doc_cluster)
```



```
1   # central points
2   print (centroids)
3   print (len(centroids))
4   print (type(centroids))
5   print (centroids.shape)
6
7   # labels
8   print (doc_cluster)
9   print (len(doc_cluster))
10  print (type(doc_cluster))
11  print (doc_cluster.shape)
12
13  labels = doc_cluster[:,0]
```

```
14  labels = list(map(int,labels))
15  print (labels[:20])
16
17  cls = set(labels)
18  print ("clusters:", cls)
```

```
1   [[ 9.807132    -1.012666     9.114005    ... -1.012666    -1.012666
2     -1.012666  ]
3    [10.46259     -0.8878248    9.421161    ... -0.8878248
    -0.8878248
4     -0.8878248 ]
5    [ 4.53212466  1.47859304  8.96488063 ... -0.49527608
    -0.56887007
6     -0.47825439]
7    ...
8    [10.77862     -1.08089      9.627346    ... -1.08089     -1.08089
9     -1.08089   ]
10   [10.02691817  8.2438794    9.13370419 ... -0.51721605
    -0.6103251
11    -0.79039602]
12   [10.45232     -0.9552607    9.6414      ... -0.9552607
    -0.9552607
13    -0.9552607 ]]
14  7
15  <class 'numpy.ndarray'>
16  (7, 5476)
17  [[2.00000000e+00 2.11885050e+04]
18   [2.00000000e+00 2.22859889e+04]
19   [2.00000000e+00 2.44449395e+04]
20   ...
21   [5.00000000e+00 4.20898259e+04]
22   [2.00000000e+00 4.61358903e+04]
23   [2.00000000e+00 4.30799672e+04]]
24  1373
25  <class 'numpy.matrix'>
26  (1373, 2)
27  [2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 5]
28  clusters: {0, 1, 2, 3, 4, 5, 6}
```

```
1   # a-4: Report top 10 words in each cluster
```

```python
import heapq

# calculate x_average
x_avg = np.mean(dw_matrix, axis=0)

vocab_df = pd.read_csv("data/science2k-vocab.txt", header=None)
vocab_array = np.array(vocab_df)
vocabs = vocab_array.tolist()
#print (vocabs[:20])
#print (len(vocabs))

def cal_cn_idx_list(n):
    i = 0
    cn = []
    cn_idx_list = []
    for p in labels:
        if p == n:
            doc = dw_matrix[i]
            cn.append(doc)
            cn_idx_list.append(i)

            if i >= len(labels)-1:
                break
        i = i + 1
    #print (cn_idx_list)
    return cn_idx_list,cn

print ("------Top 10 words in each cluster------")

for n in cls:
    cn_idx_list,cn = cal_cn_idx_list(n)
    cn_matrix = np.array(cn)
    cn_mean_matrix = np.mean(cn_matrix, axis=0)

    cn_sig_matrix = cn_mean_matrix - x_avg
    cn_sig_list = cn_sig_matrix.tolist()

    # top 10 largest numbers
    max_idx = map(cn_sig_list.index, heapq.nlargest(10,
cn_sig_list))
    max_idx_list = list(max_idx)
```

```
43        #print(list(max_idx_list))
44
45        # report top 10 words
46        top_word_list = []
47        for m in max_idx_list:
48            top_word = vocabs[m]
49            top_word_list.append(top_word)
50        print ("Cluster",n,":")
51        print (top_word_list)
```

```
1  ------Top 10 words in each cluster------
2  Cluster 0 :
3  [['eros'], ['solar'], ['kev'], ['sun'], ['elemental'],
   ['ratios'], ['asteroid'], ['ray'], ['ordinary'], ['detector']]
4  Cluster 1 :
5  [['coherence'], ['extinction'], ['patch'], ['dispersal'],
   ['coherent'], ['probabilities'], ['patches'], ['oscillations'],
   ['criteria'], ['probability']]
6  Cluster 2 :
7  [['years'], ['year'], ['scientists'], ['world'],
   ['researchers'], ['says'], ['field'], ['mail'], ['million'],
   ['focus']]
8  Cluster 3 :
9  [['co2'], ['terrestrial'], ['carbon'], ['ocean'],
   ['ecosystems'], ['atmospheric'], ['nitrogen'], ['oceans'],
   ['sink'], ['interglacial']]
10 Cluster 4 :
11 [['spectral'], ['longitude'], ['wavelengths'], ['spectra'],
   ['band'], ['wavelength'], ['incident'], ['kilometers'],
   ['elongated'], ['images']]
12 Cluster 5 :
13 [['protein'], ['cell'], ['cells'], ['expression'],
   ['proteins'], ['fig'], ['gene'], ['specific'], ['binding'],
   ['expressed']]
14 Cluster 6 :
15 [['titans'], ['clouds'], ['methane'], ['spectra'], ['cloud'],
   ['altitude'], ['atmosphere'], ['albedo'], ['flux'],
   ['saturated']]
```

```python
# a-5: Report the top ten documents that fall closest to each
cluster center

title_df = pd.read_csv("data/science2k-titles.txt",
header=None)
title_array = np.array(title_df)
titles = title_array.tolist()
#print (titles[:20])
#print (len(titles))

def get_distances(n):
    cn_idx_list,cn = cal_cn_idx_list(n)
    cn_matrix = np.array(cn)
    cn_mean_matrix = np.mean(cn_matrix, axis=0)
    distances = []
    for j in range(len(cn)):
        distances.append(euclidean_distance(cn[j],
cn_mean_matrix))
    return distances

print ("------Top 10 documents that fall closest to each
cluster center------")

for n in cls:
    cn_distances = get_distances(n)
    #print (len(cn_distances))
    #print (cn_distances)

    # top 10 smallest distances in each cluster
    min_idx = map(cn_distances.index, heapq.nsmallest(10,
cn_distances))
    min_idx_list = list(min_idx)
    #print(list(min_idx_list))

    # report top 10 documents in each cluster
    top_doc_list = []
    for m in min_idx_list:
        top_doc = titles[m]
        top_doc_list.append(top_doc)
    print ("Cluster",n,":")
    print (top_doc_list)
```

```
 1   ------Top 10 documents that fall closest to each cluster
     center------
 2   Cluster 0 :
 3   [['Archaeology in the Holy Land']]
 4   Cluster 1 :
 5   [['Archaeology in the Holy Land']]
 6   Cluster 2 :
 7   [['Similar Requirements of a Plant Symbiont and a Mammalian
     Pathogen for Prolonged Intracellular Survival'], ['A Deluge of
     Patents Creates Legal Hassles for Research'], ['Childhood
     Cancer'], ['Trojan Horses'], ['An Integrative Science Finds a
     Home'], ['An Integrative Science Finds a Home'], ['An
     Integrative Science Finds a Home'], ['Archaeology in the Holy
     Land'], ['Close Encounters: Details Veto Depth from Shadows'],
     ['Thermal, Catalytic, Regiospecific Functionalization of
     Alkanes']]
 8   Cluster 3 :
 9   [['Archaeology in the Holy Land']]
10   Cluster 4 :
11   [['Archaeology in the Holy Land'], ["Baedeker's Guide, or Just
     Plain 'Trouble'?"]]
12   Cluster 5 :
13   [['Reforming the Patent System'], ["On the Hunt for a Wolf in
     Sheep's Clothing"], ['Is Bigger Better in Cricket?'], ['Was
     Lamarck Just a Little Bit Right?'], ['Hydrogen Storage in
     Nanotubes'], ['When Pharma Merges, R&D Is the Dowry'],
     ['Superplastic Extensibility of Nanocrystalline Copper at Room
     Temperature'], ['Coupling of Stress in the ER to Activation of
     JNK Protein Kinases by Transmembrane Protein Kinase IRE1'],
     ['Mice Are Not Furry Petri Dishes'], ['<latex>$H_3^+$</latex>-
     an Ion with Many Talents']]
14   Cluster 6 :
15   [['Archaeology in the Holy Land']]
```

**a-6**

Comment on these results.

1.What has the algorithm captured?

- This algorithm mainly captures the topic information. It clusters documents, helps to get the relevant documents for each cluster center. From this, we can know which documents belong to the same topic.

2.How might such an algorithm be useful?

- This algorithm can be applied to document classification/topic clustering. By using it, we can know what topic a document belongs to and which documents belong to the same topic.
- Through the Top 10 words, we can know the key words in certain topics. Through the top 10 documents, we can find the documents that are most relevant to a certain topic.

## 1.(b)

```
1  # import modules
2
3  import numpy as np
4  import pandas as pd
5  import matplotlib.pyplot as plt
6  %matplotlib inline
7  import math
```

```
1  # load .npy file
2
3  dw_matrix = np.load('data/science2k-word-doc.npy')
4
5  print (dw_matrix)
6  print (len(dw_matrix))
7  print (type(dw_matrix))
8  print (dw_matrix.shape)
```

```
 1  [[-6.755691    -6.755691    -6.755691    ...   4.064107     5.093713
 2      3.707441   ]
 3   [-4.028205    -4.028205    -4.028205    ... -4.028205    -4.028205
 4     -4.028205   ]
 5   [-0.03370464 -1.132184    -0.03370464 ...  0.2539608    1.57568
 6      0.6594092  ]
 7   ...
 8   [-0.1301101   -0.1301101   -0.1301101   ... -0.1301101
    -0.1301101
 9     -0.1301101  ]
10   [-0.05128021 -0.05128021 -0.05128021 ... -0.05128021
    -0.05128021
```

```
11    -0.05128021]
12   [-0.06441435 -0.06441435 -0.06441435 ... -0.06441435
     -0.06441435
13    -0.06441435]]
14  5476
15  <class 'numpy.ndarray'>
16  (5476, 1373)
```

```python
1   # b-1: implement kmeans
2
3   from numpy import *
4   import matplotlib.pyplot as plt
5
6   # Compute Euclidean distance
7   def euclidean_distance(vec1, vec2):
8     eu_distance = sqrt(sum(power(vec2 - vec1, 2)))
9     return eu_distance
10
11  # init centroids with random samples
12  def init_centroids(data, k):
13    num_samples, dim = data.shape
14    centroids = zeros((k, dim))
15    for i in range(k):
16      idx = int(random.uniform(0, num_samples))
17      centroids[i, :] = data[idx, :]
18    return centroids
19
20  # kmeans
21  def kmeans(data, k):
22    num_samples = data.shape[0]
23
24    doc_cluster = mat(zeros((num_samples, 2)))
25    cls_changed = True
26
27    # init centroids
28    centroids = init_centroids(data, k)
29
30    while cls_changed:
31      cls_changed = False
32      for i in range(num_samples):
```

```
33          min_dist  = 100000.0
34          min_idx = 0
35          # for each centroid, find the centroid who is closest
36          for j in range(k):
37              distance = euclidean_distance(centroids[j, :], data[i,
    :])
38              if distance < min_dist:
39                min_dist  = distance
40                min_idx = j
41
42          # update its cluster
43          if doc_cluster[i, 0] != min_idx:
44            cls_changed = True
45            doc_cluster[i, :] = min_idx, min_dist**2
46
47        # update centroids
48        for j in range(k):
49          points_cls = data[nonzero(doc_cluster[:, 0].A == j)[0]]
50          centroids[j, :] = mean(points_cls, axis = 0)
51
52    return centroids, doc_cluster
53
54  # plot the cluster
55  def show_cluster(data, k, centroids, doc_cluster):
56    num_samples, dim = data.shape
57    if dim != 5476:
58      return 1
59
60    marks = ['or', 'ob', 'og', 'ok', '^r', '+r', 'sr', 'dr',
    '<r', 'pr','or', 'ob', 'og', 'ok', '^r', '+r', 'sr', 'dr',
    '<r', 'pr']
61    if k > len(marks):
62      return 1
63
64    # plot all samples
65    for i in range(num_samples):
66      mark_idx = int(doc_cluster[i, 0])
67      plt.plot(data[i, 0], data[i, 1], marks[mark_idx])
68
69    marks = ['Dr', 'Db', 'Dg', 'Dk', '^b', '+b', 'sb', 'db',
    '<b', 'pb','Dr', 'Db', 'Dg', 'Dk', '^b', '+b', 'sb', 'db',
    '<b', 'pb']
```
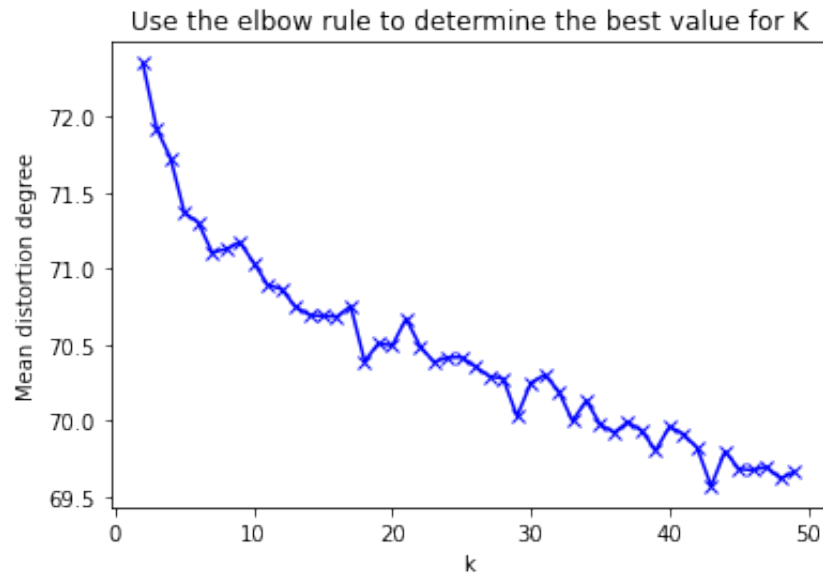
```
70    # draw the centroids
71    for i in range(k):
72       plt.plot(centroids[i, 0], centroids[i, 1], marks[i],
   markersize = 12)
73
74    plt.show()
```

```
1  # b-2: Select the best value for k
2
3  from sklearn.cluster import KMeans
4  from sklearn import metrics
5  from scipy.spatial.distance import cdist
6  import matplotlib.pyplot as plt
7
8  K = range(2, 50)
9  meandistortions = []
10
11  X = mat(dw_matrix)
12
13  for k in K:
14      kmeans = KMeans(n_clusters=k)
15      kmeans.fit(X)
16      meandistortions.append(sum(np.min(cdist(X,
   kmeans.cluster_centers_, 'euclidean'), axis=1)) / X.shape[0])
17
18  plt.plot(K, meandistortions, 'bx-')
19  plt.xlabel('k')
20  plt.ylabel(u'Mean distortion degree')
21  plt.title(u'Use the elbow rule to determine the best value for
   K');
22  plt.show()
```

Use the elbow rule to determine the best value for K



```
1   # b-3: Run the model
2   ## clustering
3   ## set k=29, because according to elbow rule, we can see that
    in the process of K value increasing, the K value corresponding
    to the position where the improvement effect of average
    distortion degree decreases the most is 29.
4
5   data = mat(dw_matrix)
6   k = 29
7   centroids, doc_cluster = kmeans(data, k)
8
9   ## plot
10  show_cluster(data, k, centroids, doc_cluster)
```

```
1   # central points
2   print (centroids)
3   print (len(centroids))
4   print (type(centroids))
5   print (centroids.shape)
6
7   # labels
8   print (doc_cluster)
9   print (len(doc_cluster))
10  print (type(doc_cluster))
11  print (doc_cluster.shape)
```

```
12
13  labels = doc_cluster[:,0]
14  labels = list(map(int,labels))
15  print (labels[:20])
16
17  cls = set(labels)
18  print ("clusters:", cls)
```

```
1   [[-0.18002037 -0.1617278  -0.11142196 ...  0.09050765
    0.19292803
2      0.18388448]
3    [-2.836106   -2.836106   -2.836106   ... -2.836106   -2.836106
4      5.681287  ]
5    [-0.62982567 -0.61657324 -0.46394096 ...  0.11452273
    -0.34449929
6     -0.38625112]
7    ...
8    [-3.18032    -3.18032    -0.96434087 ... -2.11564588
    -1.05097175
9     -2.029015  ]
10   [-1.32743722 -1.60219184 -1.32743722 ... -1.30508087
    -0.98561187
11    -0.75557164]
12   [-0.4221642  -0.4221642  -0.4221642  ... -0.4221642
    -0.4221642
13    -0.4221642 ]]
14  29
15  <class 'numpy.ndarray'>
16  (29, 1373)
17  [[2.40000000e+01 2.82920211e+04]
18   [2.60000000e+01 2.36386680e+04]
19   [0.00000000e+00 0.00000000e+00]
20   ...
21   [0.00000000e+00 0.00000000e+00]
22   [0.00000000e+00 0.00000000e+00]
23   [0.00000000e+00 0.00000000e+00]]
24  5476
25  <class 'numpy.matrix'>
26  (5476, 2)
27  [24, 26, 0, 20, 0, 26, 24, 26, 11, 0, 2, 0, 0, 24, 27, 24, 24,
    26, 4, 4]
```

```
28  clusters: {0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14,
    15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28}
```

```python
1   # b-4: Report top 10 documents in each cluster
2
3   import heapq
4
5   # calculate x_average
6   x_avg = np.mean(dw_matrix, axis=0)
7
8   title_df = pd.read_csv("data/science2k-titles.txt",
    header=None)
9   title_array = np.array(title_df)
10  titles = title_array.tolist()
11  #print (titles[:20])
12  #print (len(titles))
13
14  def cal_cn_idx_list(n):
15      i = 0
16      cn = []
17      cn_idx_list = []
18      for p in labels:
19          if p == n:
20              doc = dw_matrix[i]
21              cn.append(doc)
22              cn_idx_list.append(i)
23
24              if i >= len(labels)-1:
25                  break
26          i = i + 1
27      #print (cn_idx_list)
28      return cn_idx_list,cn
29
30  print ("------Top 10 documents------")
31
32  for n in cls:
33      cn_idx_list,cn = cal_cn_idx_list(n)
34      cn_matrix = np.array(cn)
35      cn_mean_matrix = np.mean(cn_matrix, axis=0)
36
```

```
37    cn_sig_matrix = cn_mean_matrix - x_avg
38    cn_sig_list = cn_sig_matrix.tolist()
39
40    # top 10 largest numbers
41    max_idx = map(cn_sig_list.index, heapq.nlargest(10,
   cn_sig_list))
42    max_idx_list = list(max_idx)
43    #print(list(max_idx_list))
44
45    # report top 10 documents
46    top_word_list = []
47    for m in max_idx_list:
48        top_word = titles[m]
49        top_word_list.append(top_word)
50    print ("Cluster",n,":")
51    print (top_word_list)
```

```
1  ------Top 10 documents------
2  Cluster 0 :
3  [['National Academy of Sciences Elects New Members'],
   ['Biological Control of Invading Species'], ['Scientists at
   Brookhaven'], ['Corrections and Clarifications: Timing the
   Ancestor of the HIV-1 Pandemic Strains'], ['Corrections and
   Clarifications: Timing the Ancestor of the HIV-1 Pandemic
   Strains'], ['Corrections and Clarifications: Marking Time for a
   Kingdom'], ['Corrections and Clarifications: Marking Time for a
   Kingdom'], ['Corrections and Clarifications: Marking Time for a
   Kingdom'], ['Corrections and Clarifications: One Hundred Years
   of Quantum Physics'], ['Corrections and Clarifications: One
   Hundred Years of Quantum Physics']]
4  Cluster 1 :
```

```
5   [['Global Biodiversity Scenarios for the Year 2100'],
    ['Proximity of Chromosomal Loci That Participate in Radiation-
    Induced Rearrangements in Human Cells'], ['Mate Selection and
    the Evolution of Highly Polymorphic Self/Nonself Recognition
    Genes'], ['Population Dynamical Consequences of Climate Change
    for a Small Temperate Songbird'], ['Intersubband
    Electroluminescence from Silicon-Based Quantum Cascade
    Structures'], ['On the Origin of Internal Structure of Word
    Forms'], ['Evidence for Superfluidity in Para-Hydrogen Clusters
    inside Helium-4 Droplets at 0.15 Kelvin'], ['Real-Space Imaging
    of Two-Dimensional Antiferromagnetism on the Atomic Scale'],
    ['Scanometric DNA Array Detection with Nanoparticle Probes'],
    ['The Evolutionary Fate and Consequences of Duplicate Genes']]
6   Cluster 2 :
7   [['Central Role for G Protein-Coupled Phosphoinositide 3-Kinase
    g in Inflammation'], ['Function of PI3Kg in Thymocyte
    Development, T Cell Activation, and Neutrophil Migration'],
    ['Noxa, a BH3-Only Member of the Bcl-2 Family and Candidate
    Mediator of p53-Induced Apoptosis'], ['Kinesin Superfamily
    Motor Protein KIF17 and mLin-10 in NMDA Receptor-Containing
    Vesicle Transport'], ['Requirement of JNK for Stress-Induced
    Activation of the Cytochrome c-Mediated Death Pathway'],
    ['Requirement for RORg in Thymocyte Survival and Lymphoid Organ
    Development'], ['Immune Inhibitory Receptors'], ['Role of the
    Mouse ank Gene in Control of Tissue Calcification and
    Arthritis'], ['An Oral Vaccine against NMDAR1 with Efficacy in
    Experimental Stroke and Epilepsy'], ['Ubiquitin Protein Ligase
    Activity of IAPs and Their Degradation in Proteasomes in
    Response to Apoptotic Stimuli']]
8   Cluster 3 :
9   [['DNA Damage-Induced Activation of p53 by the Checkpoint
    Kinase Chk2'], ['Timing the Radiations of Leaf Beetles:
    Hispines on Gingers from Latest Cretaceous to Recent'], ['Rapid
    Destruction of Human Cdc25A in Response to DNA Damage'],
    ['Resonant Formation of DNA Strand Breaks by Low-Energy (3 to
    20 eV) Electrons'], ['Northridge Earthquake Damage Caused by
    Geologic Focusing of Seismic Waves'], ['Radiation Tolerance of
    Complex Oxides'], ['Heightened Odds of Large Earthquakes near
    Istanbul: An Interaction-Based Probability Calculation'], ['A
    Sense of the End'], ['Stem Cells in Epithelial Tissues'],
    ['Response to RAG-Mediated V(D)J Cleavage by NBS1 and g-H2AX']]
10  Cluster 4 :
```

```
11   [['A Mouse Chronology'], ['Meltdown on Long Island'], ['Atom-
     Scale Research Gets Real'], ['Presidential Forum: Gore and Bush
     Offer Their Views on Science'], ['Help Needed to Rebuild
     Science in Yugoslavia'], ["I'd like to See America Used as a
     Global Lab"], ["Silent No Longer: 'Model Minority' Mobilizes"],
     ["Soft Money's Hard Realities"], ['Ecologists on a Mission to
     Save the World'], ['Clones: A Hard Act to Follow']]
12   Cluster 5 :
13   [['Retinal Stem Cells in the Adult Mammalian Eye'], ['Mammalian
     Neural Stem Cells'], ['From Marrow to Brain: Expression of
     Neuronal Phenotypes in Adult Mice'], ['Out of Eden: Stem Cells
     and Their Niches'], ['Turning Blood into Brain: Cells Bearing
     Neuronal Antigens Generated in Vivo from Bone Marrow'], ['The
     Genetic Program of Hematopoietic Stem Cells'], ['Genomic
     Analysis of Gene Expression in C. elegans'], ['The Initial
     Domestication of Goats (Capra hircus) in the Zagros Mountains
     10,000 Years Ago'], ['The Osteoblast: A Sophisticated
     Fibroblast under Central Surveillance'], ['Allosteric Effects
     of Pit-1 DNA Sites on Long-Term Repression in Cell Type
     Specification']]
14   Cluster 6 :
15   [['Principles for Human Gene Therapy Studies'], ['Oxidative
     Damage Linked to Neurodegeneration by Selective a-Synuclein
     Nitration in Synucleinopathy Lesions'], ['Synapses Call the
     Shots'], ['Of Chimps and Men'], ['Quantized Phonon Spectrum of
     Single-Wall Carbon Nanotubes'], ['Information Technology Takes
     a Different Tack'], ['Mothers Setting Boundaries'], ['L1
     Retrotransposons Shape the Mammalian Genome'], ['Fossils Come
     to Life in Mexico'], ['Requirement of the RNA Editing Deaminase
     ADAR1 Gene for Embryonic Erythropoiesis']]
16   Cluster 7 :
```

```
17  [['An Orientational Transition of Bent-Core Molecules in an
     Anisotropic Matrix'], ['Molecular Identification of a Taste
     Receptor Gene for Trehalose in Drosophila'], ['Multidecadal
     Changes in the Vertical Temperature Structure of the Tropical
     Troposphere'], ['Coherent High- and Low-Latitude Climate
     Variability during the Holocene Warm Period'], ['Quantum
     Criticality: Competing Ground States in Low Dimensions'],
     ['Rapid Changes in the Hydrologic Cycle of the Tropical
     Atlantic during the Last Glacial'], ['Quantum Dots as Tunable
     Kondo Impurities'], ['The Mouse House as a Recruiting Tool'],
     ['Epitopes Involved in Antibody-Mediated Protection from Ebola
     Virus'], ['Tracing the Origins of Salmonella Outbreaks']]
18  Cluster 8 :
19  [['Rapid Extragranular Plasticity in the Absence of
     Thalamocortical Plasticity in the Developing Primary Visual
     Cortex'], ['Proximity of Chromosomal Loci That Participate in
     Radiation-Induced Rearrangements in Human Cells'], ['The
     Internet of Tomorrow'], ["Detection of SO in Io's Exosphere"],
     ['The Structural Basis of Ribosome Activity in Peptide Bond
     Synthesis'], ['How Snapping Shrimp Snap: Through Cavitating
     Bubbles'], ['Learning-Induced LTP in Neocortex'], ['Heretical
     Idea Faces Its Sternest Test'], ['Cantilever Tales'], ['Single
     Photons on Demand']]
20  Cluster 9 :
21  [['How Cells Handle Cholesterol'], ['Ethanol-Induced Apoptotic
     Neurodegeneration and Fetal Alcohol Syndrome'], ['Molecular
     Evidence for the Early Evolution of Photosynthesis'], ['Direct
     Targeting of Light Signals to a Promoter Element-Bound
     Transcription Factor'], ['Architecture of RNA Polymerase II and
     Implications for the Transcription Mechanism'], ['Structural
     Evidence for Evolution of the b/a Barrel Scaffold by Gene
     Duplication and Fusion'], ['Timing the Ancestor of the HIV-1
     Pandemic Strains'], ['The Structural Basis of Ribosome Activity
     in Peptide Bond Synthesis'], ['Noxa, a BH3-Only Member of the
     Bcl-2 Family and Candidate Mediator of p53-Induced Apoptosis'],
     ['Calcium Sensitivity of Glutamate Release in a Calyx-Type
     Terminal']]
22  Cluster 10 :
```

```
23  [['Diversity and Dynamics of Dendritic Signaling'], ['Actin-
    Based Plasticity in Dendritic Spines'], ['Dopaminergic Loss and
    Inclusion Body Formation in a-Synuclein Mice: Implications for
    Neurodegenerative Disorders'], ['Untangling Dendrites with
    Quantitative Models'], ['Functional Requirement for Class I MHC
    in CNS Development and Plasticity'], ['Turning Blood into
    Brain: Cells Bearing Neuronal Antigens Generated in Vivo from
    Bone Marrow'], ['Breaking down Scientific Barriers to the Study
    of Brain and Mind'], ['Neuronal Plasticity: Increasing the Gain
    in Pain'], ['Response of Schwann Cells to Action Potentials in
    Development'], ['Kinesin Superfamily Motor Protein KIF17 and
    mLin-10 in NMDA Receptor-Containing Vesicle Transport']]
24  Cluster 11 :
25  [['Status and Improvements of Coupled General Circulation
    Models'], ['Sedimentary Rocks of Early Mars'], ['Climate
    Extremes: Observations, Modeling, and Impacts'], ["A 22,000-
    Year Record of Monsoonal Precipitation from Northern Chile's
    Atacama Desert"], ['Causes of Climate Change over the past 1000
    Years'], ['Rapid Changes in the Hydrologic Cycle of the
    Tropical Atlantic during the Last Glacial'], ['Internal
    Structure and Early Thermal Evolution of Mars from Mars Global
    Surveyor Topography and Gravity'], ['Climate Impact of Late
    Quaternary Equatorial Pacific Sea Surface Temperature
    Variations'], ['Coherent High- and Low-Latitude Climate
    Variability during the Holocene Warm Period'], ['Is El Nino
    Changing?']]
26  Cluster 12 :
27  [["Patients' Voices: The Powerful Sound in the Stem Cell
    Debate"], ['Warming of the World Ocean'], ['A Lifelong
    Fascination with the Chick Embryo'], ['Rapid Evolution of
    Reproductive Isolation in the Wild: Evidence from Introduced
    Salmon'], ['National Academy of Sciences Elects New Members'],
    ['Trans-Pacific Air Pollution'], ['Scientists at Brookhaven'],
    ['Does Science Drive the Productivity Train?'], ['Clinical
    Research'], ['Corrections and Clarifications: Luzia Is Not
    Alone']]
28  Cluster 13 :
```

29 [['Mechanism of ATP-Dependent Promoter Melting by Transcription Factor IIH'], ['b-Arrestin 2: A Receptor-Regulated MAPK Scaffold for the Activation of JNK3'], ['Role for Rapid Dendritic Protein Synthesis in Hippocampal mGluR-Dependent Long-Term Depression'], ['Interacting Molecular Loops in the Mammalian Circadian Clock'], ["Packard Heir Signs up for National 'Math Wars'"], ['Virus-Induced Neuronal Apoptosis Blocked by the Herpes Simplex Virus Latency-Associated Transcript'], ['Transgenic Mouse Model of Stunned Myocardium'], ['Interconnected Feedback Loops in the Neurospora Circadian System'], ['Inhibition of Adipogenesis by Wnt Signaling'], ['An Inherited Functional Circadian Clock in Zebrafish Embryos']]
30 Cluster 14 :
31 [['N-Cadherin, a Cell Adhesion Molecule Involved in Establishment of Embryonic Left-Right Asymmetry'], ['Transgenic Mouse Model of Stunned Myocardium'], ['Nota Bene: Contortions of the Heart'], ['Cardiovascular Evidence for an Intermediate or Higher Metabolic Rate in an Ornithischian Dinosaur'], ['Resetting of Circadian Time in Peripheral Tissues by Glucocorticoid Signaling'], ['NIH, under Pressure, Boosts Minority Health Research'], ['Generalized Potential of Adult Neural Stem Cells'], ['New Age Semiconductors Pick up the Pace'], ["Stress: The Invisible Hand in Eastern Europe's Death Rates"], ['Role of Adenine Nucleotide Translocator 1 in mtDNA Maintenance']]
32 Cluster 15 :
33 [['Emerging Infectious Diseases of Wildlife-Threats to Biodiversity and Human Health'], ['A Tale of Two Futures: HIV and Antiretroviral Therapy in San Francisco'], ['Predictions of Biodiversity Response to Genetically Modified Herbicide-Tolerant Crops'], ['A Tale of Two Selves'], ['Origins of HIV'], ['Fairness versus Reason in the Ultimatum Game'], ['One Sequence, Two Ribozymes: Implications for the Emergence of New Ribozyme Folds'], ['Reversal of Antipsychotic-Induced Working Memory Deficits by Short-Term Dopamine D1 Receptor Stimulation'], ['Potent Analgesic Effects of GDNF in Neuropathic Pain States'], ['AIDS in a New Millennium']]
34 Cluster 16 :

35  [["Evidence for Crystalline Water and Ammonia Ices on Pluto's Satellite Charon"], ['Mount St. Helens, Master Teacher'], ['Atomic Layer Deposition of Oxide Thin Films with Metal Alkoxides as Oxygen Sources'], ['Dynamics of the Pacific-North American Plate Boundary in the Western United States'], ['Discovery of a Transient Absorption Edge in the X-ray Spectrum of GRB 990705'], ['Differential Clustering of CD4 and CD3z during T Cell Recognition'], ['Memory-A Century of Consolidation'], ['Motility Powered by Supramolecular Springs and Ratchets'], ['Response of Schwann Cells to Action Potentials in Development'], ['Piecing Together the Biggest Puzzle of All']]

36  Cluster 17 :

37  [['An Arresting Start for MAPK'], ["CERN's Gamble Shows Perils, Rewards of Playing the Odds"], ["Outrageous Events: Don't Count Them out"], ['Mechanism of ATP-Dependent Promoter Melting by Transcription Factor IIH'], ['A Wetter, Younger Mars Emerging'], ['Rapid Extragranular Plasticity in the Absence of Thalamocortical Plasticity in the Developing Primary Visual Cortex'], ['Rounding out Solutions to Three Conjectures'], ['On the Origin of Internal Structure of Word Forms'], ['Alternative Views on Alternative Medicine'], ["Candida's Arranged Marriage"]]

38  Cluster 18 :

39  [['Mode-Specific Energy Disposal in the Four-Atom Reaction <latex>$OH + D_2 \\rightarrow HOD + D$</latex>'], ['A Short Fe-Fe Distance in Peroxodiferric Ferritin: Control of Fe Substrate versus Cofactor Decay?'], ['The Evolutionary Fate and Consequences of Duplicate Genes'], ["Fermat's Last Theorem's First Cousin"], ['Social Mentalizing Abilities in Mental Patients'], ['Nonavian Feathers in a Late Triassic Archosaur'], ['Nonbiological Fractionation of Iron Isotopes'], ['A Cyclic Carbanionic Valence Isomer of a Carbocation: Diphosphino Analogs of Diaminocarbocations'], ['Mechanisms of Ordering in Striped Patterns'], ['Bioinformatics in the Information Age']]

40  Cluster 19 :

```
41   [['Uptake of Glutamate into Synaptic Vesicles by an Inorganic
     Phosphate Transporter'], ['Gatekeepers of the Nucleus'],
     ['VirB/D4-Dependent Protein Translocation from Agrobacterium
     into Plant Cells'], ['Structure of the Light-Driven Chloride
     Pump Halorhodopsin at 1.8 <latex>$\\AA$</latex> Resolution'],
     ['Rab1 Recruitment of p115 into a cis-SNARE Complex:
     Programming Budding COPII Vesicles for Fusion'], ['Connectivity
     of Marine Populations: Open or Closed?'], ['How Cells Handle
     Cholesterol'], ['Trans-Pacific Air Pollution'], ['Transmembrane
     Molecular Pump Activity of Niemann-Pick C1 Protein'], ['The
     Influence of Canadian Forest Fires on Pollutant Concentrations
     in the United States']]
42   Cluster 20 :
43   [['Crystal Structure of the Ribonucleoprotein Core of the
     Signal Recognition Particle'], ['The Complete Atomic Structure
     of the Large Ribosomal Subunit at 2.4 <latex>$\\AA$</latex>
     Resolution'], ['Three-Dimensional Structure of the Tn5 Synaptic
     Complex Transposition Intermediate'], ['The Structural Basis of
     Ribosome Activity in Peptide Bond Synthesis'], ['Architecture
     of RNA Polymerase II and Implications for the Transcription
     Mechanism'], ['Comparative Genomics of the Eukaryotes'],
     ['Positional Syntenic Cloning and Functional Characterization
     of the Mammalian Circadian Mutation tau'], ['The Way Things
     Move: Looking under the Hood of Molecular Motor Proteins'],
     ['The Genome Sequence of Drosophila melanogaster'], ['Structure
     of the RNA Polymerase Domain of E. coli Primase']]
44   Cluster 21 :
45   [['ORCA3, a Jasmonate-Responsive Transcriptional Regulator of
     Plant Primary and Secondary Metabolism'], ['Psychological and
     Neural Mechanisms of the Affective Dimension of Pain'], ['The
     Complete Atomic Structure of the Large Ribosomal Subunit at 2.4
     <latex>$\\AA$</latex> Resolution'], ['Green, Catalytic
     Oxidation of Alcohols in Water'], ['One Sequence, Two
     Ribozymes: Implications for the Emergence of New Ribozyme
     Folds'], ['A Structural Model of Transcription Elongation'],
     ['Template Boundary in a Yeast Telomerase Specified by RNA
     Structure'], ['Impacts of Climatic Change and Fishing on
     Pacific Salmon Abundance over the past 300 Years'], ['Evidence
     for Recent Groundwater Seepage and Surface Runoff on Mars'],
     ['Calcium-Aluminum-Rich Inclusions from Enstatite Chondrites:
     Indigenous or Foreign?']]
46   Cluster 22 :
```

```
47   [["Packard Heir Signs up for National 'Math Wars'"], ['Islamic
     Women in Science'], ['Not (Just) in Kansas Anymore'],
     ['Graduate Educators Struggle to Grade Themselves'], ['The
     Spirit of Discovery'], ['Support Grows for British Exercise to
     Allocate University Funds'], ['Sharp Jump in Teaching Fellows
     Draws Fire from Educators'], ['Presidential Forum: Gore and
     Bush Offer Their Views on Science'], ["Iran's Scientists
     Cautiously Reach out to the World"], ['Scaling up HIV/AIDS
     Programs to National Coverage']]
48   Cluster 23 :
49   [['Advances in the Physics of High-Temperature
     Superconductivity'], ['Quantum Criticality: Competing Ground
     States in Low Dimensions'], ['Orbital Physics in Transition-
     Metal Oxides'], ['The Atom-Cavity Microscope: Single Atoms
     Bound in Orbit by Single Photons'], ["Negative Poisson's Ratios
     for Extreme States of Matter"], ['Self-Mode-Locking of Quantum
     Cascade Lasers with Giant Ultrafast Optical Nonlinearities'],
     ['Generating Solitons by Phase Engineering of a Bose-Einstein
     Condensate'], ['Imaging Precessional Motion of the
     Magnetization Vector'], ['Subatomic Features on the Silicon
     (111)-(7 x 7) Surface Observed by Atomic Force Microscopy'],
     ['Blue-Fluorescent Antibodies']]
50   Cluster 24 :
51   [['NEAR at Eros: Imaging and Spectral Results'], ['Climate
     Extremes: Observations, Modeling, and Impacts'], ['Reduction of
     Tropical Cloudiness by Soot'], ['Causes of Climate Change over
     the past 1000 Years'], ['The Atom-Cavity Microscope: Single
     Atoms Bound in Orbit by Single Photons'], ['High Magma Storage
     Rates before the 1983 Eruption of Kilauea, Hawaii'],
     ['Hematopoietic Stem Cell Quiescence Maintained by
     <latex>$p21^{cip1/waf1}$</latex>'], ['Climate Impact of Late
     Quaternary Equatorial Pacific Sea Surface Temperature
     Variations'], ['Internal Structure and Early Thermal Evolution
     of Mars from Mars Global Surveyor Topography and Gravity'],
     ['Isotope Fractionation and Atmospheric Oxygen: Implications
     for Phanerozoic <latex>$O_2$</latex> Evolution']]
52   Cluster 25 :
```

53   [['Can Protected Areas Be Expanded in Africa?'], ['Status and
     Improvements of Coupled General Circulation Models'],
     ['Surveying the SBIR Program'], ['Tissue Engineers Build New
     Bone'], ['Interfering with Gene Expression'], ['Corrections and
     Clarifications: Commercialization of Genetic Research and
     Public Policy'], ['Corrections and Clarifications:
     Commercialization of Genetic Research and Public Policy'],
     ['Thyroid Tumor Banks'], ['Corrections and Clarifications: The
     Global Spread of Malaria in a Future, Warmer World'],
     ["Corrections and Clarifications: Fermat's Last Theorem's First
     Cousin"]]
54   Cluster 26 :
55   [['Global Analysis of the Genetic Network Controlling a
     Bacterial Cell Cycle'], ['Mitotic Misregulation and Human
     Aging'], ['Genes Expressed in Human Tumor Endothelium'],
     ['Inhibition of Adipogenesis by Wnt Signaling'], ['Interacting
     Molecular Loops in the Mammalian Circadian Clock'], ['From
     Marrow to Brain: Expression of Neuronal Phenotypes in Adult
     Mice'], ['Noxa, a BH3-Only Member of the Bcl-2 Family and
     Candidate Mediator of p53-Induced Apoptosis'], ['Translocation
     of C. elegans CED-4 to Nuclear Membranes during Programmed Cell
     Death'], ['Stat3-Mediated Transformation of NIH-3T3 Cells by
     the Constitutively Active Q205L <latex>$G\\alpha_o$</latex>
     Protein'], ['A Subset of Viral Transcripts Packaged within
     Human Cytomegalovirus Particles']]
56   Cluster 27 :
57   [['Positional Syntenic Cloning and Functional Characterization
     of the Mammalian Circadian Mutation tau'], ['The Genome
     Sequence of Drosophila melanogaster'], ['Gridlock, an HLH Gene
     Required for Assembly of the Aorta in Zebrafish'], ['Pif1p
     Helicase, a Catalytic Inhibitor of Telomerase in Yeast'],
     ['Conservation and Novelty in the Evolution of Cell Adhesion
     and Extracellular Matrix Genes'], ['Requirement of Mis6
     Centromere Connector for Localizing a CENP-A-Like Protein in
     Fission Yeast'], ['Role of the Mouse ank Gene in Control of
     Tissue Calcification and Arthritis'], ['Comparative Genomics of
     the Eukaryotes'], ['Resetting of Circadian Time in Peripheral
     Tissues by Glucocorticoid Signaling'], ['Accumulation of
     Dietary Cholesterol in Sitosterolemia Caused by Mutations in
     Adjacent ABC Transporters']]
58   Cluster 28 :

```
59  [['New Observational Constraints for Atmospheric Hydroxyl on
     Global and Hemispheric Scales'], ['Dimensionality Effects in
     the Lifetime of Surface States'], ['Forster Energy Transfer in
     an Optical Microcavity'], ['Quantitative Imaging of Lateral
     ErbB1 Receptor Signal Propagation in the Plasma Membrane'], ['A
     Potent Greenhouse Gas Identified in the Atmosphere:
     $SF_5CF_3$'], ['Why Stem Cells?'], ['Intersubband
     Electroluminescence from Silicon-Based Quantum Cascade
     Structures'], ['Blue-Fluorescent Antibodies'], ['Molecules in a
     Bose-Einstein Condensate'], ['Quantifying Denitrification and
     Its Effect on Ozone Recovery']]
```

```python
1   # b-5: Report the top ten words
2
3   vocab_df = pd.read_csv("data/science2k-vocab.txt", header=None)
4   vocab_array = np.array(vocab_df)
5   vocabs = vocab_array.tolist()
6   #print (vocabs[:20])
7   #print (len(vocabs))
8
9   def get_distances(n):
10      cn_idx_list,cn = cal_cn_idx_list(n)
11      cn_matrix = np.array(cn)
12      cn_mean_matrix = np.mean(cn_matrix, axis=0)
13      distances = []
14      for j in range(len(cn)):
15          distances.append(euclidean_distance(cn[j],
    cn_mean_matrix))
16      return distances
17
18  print ("------Top 10 words------")
19
20  for n in cls:
21      cn_distances = get_distances(n)
22      #print (len(cn_distances))
23      #print (cn_distances)
24
25      # top 10 smallest distances in each cluster
26      min_idx = map(cn_distances.index, heapq.nsmallest(10,
    cn_distances))
```

```
27      min_idx_list = list(min_idx)
28      #print(list(min_idx_list))
29
30      # report top 10 words
31      top_doc_list = []
32      for m in min_idx_list:
33          top_doc = vocabs[m]
34          top_doc_list.append(top_doc)
35      print ("Cluster",n,":")
36      print (top_doc_list)
```

```
1   ------Top 10 words------
2   Cluster 0 :
3   [['graduate'], ['cyclin'], ['historical'], ['magnification'],
    ['switching'], ['radial'], ['apart'], ['digital'],
    ['stronger'], ['tuning']]
4   Cluster 1 :
5   [['fig']]
6   Cluster 2 :
7   [['amount'], ['identical'], ['patterns'], ['cortex'],
    ['origin'], ['correspondence'], ['estimates'], ['working'],
    ['thin'], ['voltage']]
8   Cluster 3 :
9   [['fig']]
10  Cluster 4 :
11  [['medium'], ['june'], ['spin'], ['united'], ['cellular'],
    ['century'], ['method'], ['step'], ['helix'], ['increasing']]
12  Cluster 5 :
13  [['fig'], ['fig']]
14  Cluster 6 :
15  [['fig'], ['fig']]
16  Cluster 7 :
17  [['fig']]
18  Cluster 8 :
19  [['fig']]
20  Cluster 9 :
21  [['start'], ['vol'], ['www'], ['end'], ['cells'], ['time'],
    ['data'], ['cell'], ['two'], ['science']]
22  Cluster 10 :
23  [['cell'], ['data'], ['time'], ['science'], ['two'],
    ['protein'], ['end'], ['fig'], ['cells']]
24  Cluster 11 :
```

```
25  [['view'], ['band'], ['developed'], ['atoms'], ['chemical'],
    ['crystal'], ['correspondence'], ['expressed'], ['isolated'],
    ['oxygen']]
26  Cluster 12 :
27  [['fig']]
28  Cluster 13 :
29  [['fig']]
30  Cluster 14 :
31  [['fig']]
32  Cluster 15 :
33  [['fig']]
34  Cluster 16 :
35  [['fig']]
36  Cluster 17 :
37  [['fig']]
38  Cluster 18 :
39  [['fig']]
40  Cluster 19 :
41  [['fig']]
42  Cluster 20 :
43  [['information'], ['sequences'], ['neurons'], ['point'],
    ['determined'], ['age'], ['current'], ['addition'],
    ['activation'], ['expressed']]
44  Cluster 21 :
45  [['fig']]
46  Cluster 22 :
47  [['fig']]
48  Cluster 23 :
49  [['metal'], ['black'], ['chemical'], ['atoms'], ['involved'],
    ['open'], ['transition'], ['class'], ['components'],
    ['measurements']]
50  Cluster 24 :
51  [['science'], ['nature'], ['change'], ['site'], ['specific'],
    ['studies'], ['molecular'], ['changes'], ['says'], ['long']]
52  Cluster 25 :
53  [['fig']]
54  Cluster 26 :
55  [['science'], ['end'], ['cell'], ['cells'], ['fig'],
    ['protein'], ['data'], ['two']]
56  Cluster 27 :
```

```
57  [['results'], ['university'], ['shown'], ['genes'],
    ['observed'], ['different'], ['type'], ['expression'],
    ['high'], ['dna']]
58  Cluster 28 :
59  [['fig']]
```

**b-6**

Comment on these results.

1.How might such an algorithm be useful?

- This algorithm can be applied to subject retrieval/document retrieval, which is to retrieve documents based on keywords. For example, if I want to search for science fiction documents, I can find the relevant documents by searching for the keyword "science fiction".

2.What is different about clustering terms from clustering documents?

- The difference is that clustering terms can retrieve documents by keywords, but clustering documents cannot retrieve documents by keywords.

# Question 2

```python
1  import numpy as np
2  import matplotlib.pyplot as plt
3  import math
4  import random
```

**2.(a) Download the Old Faithful Geyser Dataset. The data file contains 272 observations of (eruption time, waiting time). Treat each entry as a 2 dimensional feature vector. Parse and plot all data points on 2-D plane.**

```python
1  Data_list = []
2  with open("old.txt", 'r') as in_file:
3      for line in in_file.readlines():
4          point = []
5          point.append(float(line.split()[1]))
```

```
 6                point.append(float(line.split()[2]))
 7                Data_list.append(point)
 8    Data = np.array(Data_list)
 9    eruptions=[]
10    waiting=[]
11    for i in range(len(Data)):
12        eruptions.append(Data[i][0])
13        waiting.append(Data[i][1])
14
15    #plot scatter
16    plt.title('Old Faithful Geyser')
17    plt.scatter(eruptions, waiting)
18    plt.xlabel('eruptions')
19    plt.ylabel('waiting')
20    plt.show()
```



Old Faithful Geyser

As can be seen from the scatter plot, the data set has two categories.

**2.(b) Implement a bimodal GMM model to fit all data points using EM algorithm. Explain the reasoning behind your termination criteria. For this problem, we assume the covariance matrix is spherical (i.e., it has the form of σ^2I for scalar σ) and you can randomly initialize Gaussian parameters. For evaluation purposes, please submit the following figures:**

Since the initial guess of mu and sigma are whole numbers, I set the termination criteria=0.01. When the difference between the two iterations is less than 0.01, it is considered that the parameters no longer change.

```python
parameter_dict = {}
parameter_dict["mu1"] = np.array([0, 0])
parameter_dict["sigma1"] = np.array([[1, 0], [0, 1]])
parameter_dict["mu2"] = np.array([0, 0])
parameter_dict["sigma2"] = np.array([[1, 0], [0, 1]])
parameter_dict["piweight"] = 0.5
parameter_dict["gama_list"] = []


def set_parameter(mu_1, sigma_1, mu_2, sigma_2, pi_weight):
    parameter_dict["mu1"] = mu_1
    parameter_dict["mu1"].shape = (2, 1)
    parameter_dict["sigma1"] = sigma_1
    parameter_dict["mu2"] = mu_2
    parameter_dict["mu2"].shape = (2, 1)
    parameter_dict["sigma2"] = sigma_2
    parameter_dict["piweight"] = pi_weight


def PDF(data, mu, sigma):

    sigma_sqrt = math.sqrt(np.linalg.det(sigma))
    sigma_inv = np.linalg.inv(sigma)
    data.shape = (2, 1)
    mu.shape = (2, 1)
    minus_mu = data - mu
    minus_mu_trans = np.transpose(minus_mu)
    res = (1.0 / (2.0 * math.pi * sigma_sqrt)) * math.exp(
        (-0.5) * (np.dot(np.dot(minus_mu_trans, sigma_inv),
    minus_mu)))
    return res

def E_step(Data):
    sigma_1 = parameter_dict["sigma1"]
    sigma_2 = parameter_dict["sigma2"]
    pw = parameter_dict["piweight"]
    mu_1 = parameter_dict["mu1"]
    mu_2 = parameter_dict["mu2"]
```

```python
38
39    parameter_dict["gama_list"] = []
40    for point in Data:
41        gama_i = (pw * PDF(point, mu_2, sigma_2)) / (
42            (1.0 - pw) * PDF(point, mu_1, sigma_1) + pw *
   PDF(point, mu_2, sigma_2))
43        parameter_dict["gama_list"].append(gama_i)
44
45
46 def M_step(Data):
47    N1 = 0
48    N2 = 0
49    for i in range(len(parameter_dict["gama_list"])):
50        N1 += 1.0 - parameter_dict["gama_list"][i]
51        N2 += parameter_dict["gama_list"][i]
52
53    new_mu_1 = np.array([0, 0])
54    new_mu_2 = np.array([0, 0])
55    for i in range(len(parameter_dict["gama_list"])):
56        new_mu_1 = new_mu_1 + Data[i] * (1 -
   parameter_dict["gama_list"][i]) / N1
57        new_mu_2 = new_mu_2 + Data[i] *
   parameter_dict["gama_list"][i] / N2
58
59
60    new_mu_1.shape = (2, 1)
61    new_mu_2.shape = (2, 1)
62
63    new_sigma_1 = np.array([[0, 0], [0, 0]])
64    new_sigma_2 = np.array([[0, 0], [0, 0]])
65    for i in range(len(parameter_dict["gama_list"])):
66        data_tmp = [0, 0]
67        data_tmp[0] = Data[i][0]
68        data_tmp[1] = Data[i][1]
69        vec_tmp = np.array(data_tmp)
70        vec_tmp.shape = (2, 1)
71        new_sigma_1 = new_sigma_1 + np.dot((vec_tmp -
   new_mu_1), (vec_tmp - new_mu_1).transpose()) * (1.0 -
   parameter_dict["gama_list"][i]) / N1
72        new_sigma_2 = new_sigma_2 + np.dot((vec_tmp -
   new_mu_2), (vec_tmp - new_mu_2).transpose()) *
   parameter_dict["gama_list"][i] / N2
```

```
73        # print np.dot((vec_tmp-new_mu_1), (vec_tmp-
   new_mu_1).transpose())
74     new_pi = N2 / len(parameter_dict["gama_list"])
75
76     parameter_dict["mu1"] = new_mu_1
77     parameter_dict["mu2"] = new_mu_2
78     parameter_dict["sigma1"] = new_sigma_1
79     parameter_dict["sigma2"] = new_sigma_2
80     parameter_dict["piweight"] = new_pi
81
```

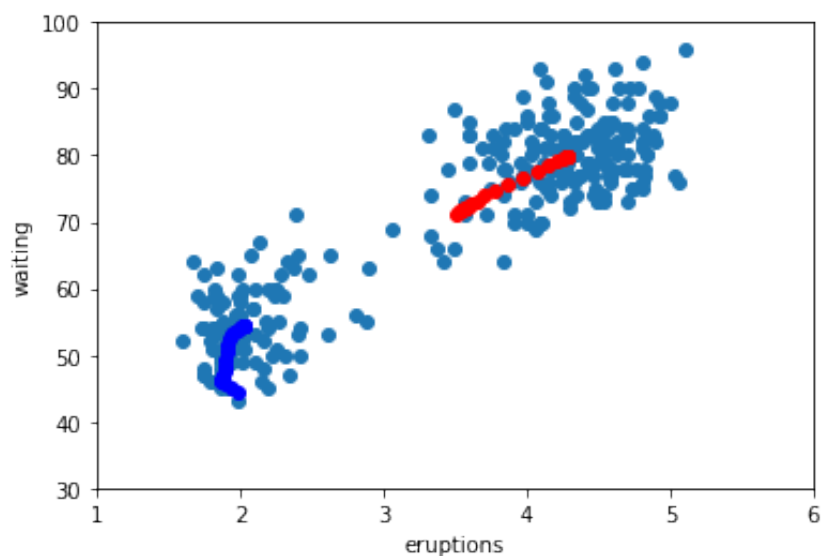## i). Plot the trajectories of two mean vectors in 2 dimensions (i.e., coordinates vs. iteration).

```
1  #Plot the trajectories
2  def EM_iterate_trajectories(iter_time, Data, mu_1, sigma_1,
   mu_2, sigma_2, pi_weight, esp=0.01):
3
4      mean_trace_1 = [[], []]
5      mean_trace_2 = [[], []]
6
7      set_parameter(mu_1, sigma_1, mu_2, sigma_2, pi_weight)
8      if iter_time == None:
9          while (True):
10             old_mu_1 = parameter_dict["mu1"].copy()
11             old_mu_2 = parameter_dict["mu2"].copy()
12             E_step(Data)
13             M_step(Data)
14             delta_1 = parameter_dict["mu1"] - old_mu_1
15             delta_2 = parameter_dict["mu2"] - old_mu_2
16
17             mean_trace_1[0].append(parameter_dict["mu1"][0][0])
18             mean_trace_1[1].append(parameter_dict["mu1"][1][0])
19             mean_trace_2[0].append(parameter_dict["mu2"][0][0])
20             mean_trace_2[1].append(parameter_dict["mu2"][1][0])
21             if math.fabs(delta_1[0]) < esp and
   math.fabs(delta_1[1]) < esp and math.fabs(
22                     delta_2[0]) < esp and math.fabs(delta_2[1])
   < esp:
23                 break
24     else:
```

```
25              for i in range(iter_time):
26                  pass
27
28          plt.xlim(xmax=6, xmin=1)
29          plt.ylim(ymax=100, ymin=30)
30          plt.xlabel("eruptions")
31          plt.ylabel("waiting")
32          plt.scatter(eruptions, waiting)
33          plt.plot(mean_trace_1[0], mean_trace_1[1], 'ro')
34          plt.plot(mean_trace_2[0], mean_trace_2[1], 'bo')
35          plt.show()
36
37   def task_1():
38          Mu_1 = np.array([3, 60])
39          Sigma_1 = np.array([[10, 0], [0, 10]])
40          Mu_2 = np.array([1, 30])
41          Sigma_2 = np.array([[10, 0], [0, 10]])
42          Pi_weight = 0.5
43          EM_iterate_trajectories(None, Data, Mu_1, Sigma_1, Mu_2,
     Sigma_2, Pi_weight)
44
45   task_1()
```



Choose $\mu_1 = [3, 60], \mu_2 = [1, 30], \sigma_1 = [[10, 0], [0, 10]], \sigma_2 = [[10, 0], [0, 10]]$ as the initial input since after observing the scatter plot of the dataset we can see the two $\mu$ we guess is on the left side of the bottom of two categories. Thus, we could see the trajectories clearly and as complete as possible.

## ii). Run your program for 50 times with different initial parameter guesses. Show the distribution of the total number of iterations needed for algorithm to converge.

```python
#Run the program for 50 times with different initial parameter
guesses.
def EM_iterate_times(Data, mu_1, sigma_1, mu_2, sigma_2,
pi_weight, esp=0.01):

    set_parameter(mu_1, sigma_1, mu_2, sigma_2, pi_weight)
    iter_times = 0
    while (True):
        iter_times += 1
        old_mu_1 = parameter_dict["mu1"].copy()
        old_mu_2 = parameter_dict["mu2"].copy()
        E_step(Data)
        M_step(Data)
        delta_1 = parameter_dict["mu1"] - old_mu_1
        delta_2 = parameter_dict["mu2"] - old_mu_2
        if math.fabs(delta_1[0]) < esp and
math.fabs(delta_1[1]) < esp and math.fabs(
                delta_2[0]) < esp and math.fabs(delta_2[1]) <
esp:
            break
    return iter_times

def task_2():
    iter_times=0
    try:
        x_11 = random.uniform(0, 4)
        x_12 = random.uniform(30, 70)
        x_21 = random.uniform(2, 6)
        x_22 = random.uniform(50, 100)
        Mu_1 = np.array([x_11, x_12])
        x_31 = random.uniform(1, 20)
        Sigma_1 = np.array([[x_31, 0], [0, x_31]])
        Mu_2 = np.array([x_21, x_22])
        Sigma_2 = np.array([[x_31, 0], [0, x_31]])
        Pi_weight = 0.5
```

```python
32          iter_times = EM_iterate_times(Data, Mu_1, Sigma_1,
    Mu_2, Sigma_2, Pi_weight)
33          print ('mu1=[',x_11,',',x_12,'],mu2=
    [',x_21,',',x_22,']','total number of iterations=',iter_times)
34      except Exception:
35          print (Exception)
36      return iter_times
37
38  iteration=[]
39  for i in range(50):
40      iteration.append(task_2())
41  plt.hist(iteration)
42  plt.title('GMM')
43  plt.show
44
45
46
47  def M_step_1(Data):
48      N1 = 0
49      N2 = 0
50      for i in range(len(parameter_dict["gama_list"])):
51          N1 += 1.0 - parameter_dict["gama_list"][i]
52          N2 += parameter_dict["gama_list"][i]
53
54      new_mu_1 = np.array([0, 0])
55      new_mu_2 = np.array([0, 0])
56      for i in range(len(parameter_dict["gama_list"])):
57          new_mu_1 = new_mu_1 + Data[i] * (1 -
    parameter_dict["gama_list"][i]) / N1
58          new_mu_2 = new_mu_2 + Data[i] *
    parameter_dict["gama_list"][i] / N2
59
60
61      new_mu_1.shape = (2, 1)
62      new_mu_2.shape = (2, 1)
63
64      new_sigma_1 = np.array([[0, 0], [0, 0]])
65      new_sigma_2 = np.array([[0, 0], [0, 0]])
66      for i in range(len(parameter_dict["gama_list"])):
67          data_tmp = [0, 0]
68          data_tmp[0] = Data[i][0]
69          data_tmp[1] = Data[i][1]
```

```
70          vec_tmp = np.array(data_tmp)
71          vec_tmp.shape = (2, 1)
72          new_sigma_1 = new_sigma_1 + np.dot((vec_tmp -
   new_mu_1), (vec_tmp - new_mu_1).transpose()) * (1.0 -
   parameter_dict["gama_list"][i]) / N1
73          new_sigma_2 = new_sigma_2 + np.dot((vec_tmp -
   new_mu_2), (vec_tmp - new_mu_2).transpose()) *
   parameter_dict["gama_list"][i] / N2
74          # print np.dot((vec_tmp-new_mu_1), (vec_tmp-
   new_mu_1).transpose())
75      new_pi = N2 / len(parameter_dict["gama_list"])
76
77      parameter_dict["mu1"] = new_mu_1
78      parameter_dict["mu2"] = new_mu_2
79      parameter_dict["sigma1"] = new_sigma_1
80      parameter_dict["sigma2"] = new_sigma_2
81      parameter_dict["piweight"] = new_pi
82      covm1=[]
83      covm2=[]
84      print(new_mu_1)
85      print(new_mu_2)
86      covm1.append(new_sigma_1[0][0])
87      covm1.append(new_sigma_1[1][1])
88      covm2.append(new_sigma_2[0][0])
89      covm2.append(new_sigma_2[1][1])
90      print(covm1)
91      print(covm2)
92  #EM result:center and covariance matrix
93  E_step(Data)
94  M_step_1(Data)
```

```
1  mu1=[ 0.1893135405201245 , 35.369020344606255 ],mu2=[
   4.345307785548732 , 64.63663783179955 ] total number of
   iterations= 15
2  mu1=[ 2.221758431528942 , 45.47747359129394 ],mu2=[
   4.0934557777931495 , 53.714273021455675 ] total number of
   iterations= 15
3  mu1=[ 2.8432992305101483 , 31.42269914464936 ],mu2=[
   4.31931126754017 , 83.83778519912337 ] total number of
   iterations= 10
```

```
 4  mu1=[ 3.094910237592582 , 53.37228667775115 ],mu2=[
    4.619537418786634 , 64.63983267902262 ] total number of
    iterations= 9
 5  mu1=[ 0.8719747538166236 , 46.76388076310295 ],mu2=[
    2.005938861877005 , 91.38603886689017 ] total number of
    iterations= 6
 6  mu1=[ 0.0845258727660041 , 59.246730062985264 ],mu2=[
    2.49894764043518 , 60.130329127166405 ] total number of
    iterations= 13
 7  mu1=[ 2.9629717440087098 , 53.06543041297001 ],mu2=[
    4.555527171270407 , 86.03855911794619 ] total number of
    iterations= 6
 8  mu1=[ 1.9769560409458733 , 44.13450566893186 ],mu2=[
    3.041613249704096 , 77.49528397696082 ] total number of
    iterations= 8
 9  mu1=[ 1.9968414447742209 , 31.062437605032972 ],mu2=[
    5.274612911990247 , 55.565061254341444 ] total number of
    iterations= 9
10  mu1=[ 2.2803772969098466 , 33.976268025059234 ],mu2=[
    3.8430136950093567 , 76.18116984311325 ] total number of
    iterations= 11
11  mu1=[ 1.1376652390642539 , 45.43551455238631 ],mu2=[
    4.36609105083032 , 84.43625665492587 ] total number of
    iterations= 3
12  mu1=[ 3.6629922126674384 , 46.86452042143449 ],mu2=[
    2.2748867821848564 , 52.294834501504766 ] total number of
    iterations= 16
13  mu1=[ 0.3915060586649006 , 69.2662185623231 ],mu2=[
    3.268545013895078 , 88.83529026087001 ] total number of
    iterations= 13
14  mu1=[ 2.630806312503527 , 39.35272269523833 ],mu2=[
    5.220488801727269 , 52.28495969819305 ] total number of
    iterations= 17
15  mu1=[ 1.7861812550507268 , 35.90250627140895 ],mu2=[
    3.9633365219410694 , 94.03209006667464 ] total number of
    iterations= 4
16  mu1=[ 0.3691169903609919 , 48.44419814390284 ],mu2=[
    5.229180814491622 , 70.2663021636601 ] total number of
    iterations= 8
17  mu1=[ 0.8603210513259905 , 46.90288032657365 ],mu2=[
    4.233308171087751 , 51.212466118429504 ] total number of
    iterations= 16
```

```
18  mu1=[ 0.4911690066962877 , 66.02275700646787 ],mu2=[
    2.944335115977025 , 97.9389785646695 ] total number of
    iterations= 16
19  mu1=[ 0.8051782769011528 , 48.5732584441359 ],mu2=[
    2.546086943527769 , 93.16077946406496 ] total number of
    iterations= 7
20  mu1=[ 3.4854843931704442 , 47.73486188316144 ],mu2=[
    5.979177585377549 , 61.28617570938973 ] total number of
    iterations= 11
21  mu1=[ 2.852386438452337 , 36.637774289861255 ],mu2=[
    4.383406686700793 , 71.54195119524215 ] total number of
    iterations= 12
22  mu1=[ 2.0043255070479726 , 36.400725821714545 ],mu2=[
    5.739160437971196 , 73.44246797647395 ] total number of
    iterations= 11
23  mu1=[ 2.6957374475102016 , 31.539154708990512 ],mu2=[
    5.232567282361543 , 84.33006323297006 ] total number of
    iterations= 10
24  mu1=[ 3.3878592670906844 , 38.89144997455178 ],mu2=[
    4.7804473273958 , 63.991090254053475 ] total number of
    iterations= 14
25  mu1=[ 3.513569268511476 , 31.790252280542767 ],mu2=[
    5.613092774672461 , 68.48760949542064 ] total number of
    iterations= 15
26  mu1=[ 0.9846761978958227 , 45.85074014878293 ],mu2=[
    5.893955561355381 , 88.01744812757072 ] total number of
    iterations= 5
27  mu1=[ 1.6084893867380616 , 69.44175084065354 ],mu2=[
    2.099294439324174 , 82.90147379022544 ] total number of
    iterations= 10
28  mu1=[ 2.572486164344156 , 41.12062253180023 ],mu2=[
    2.226112724849305 , 93.70290690390274 ] total number of
    iterations= 5
29  mu1=[ 0.9535833498412507 , 57.176492968114694 ],mu2=[
    5.568121493274575 , 54.1627517305873 ] total number of
    iterations= 11
30  mu1=[ 3.318145031187072 , 42.376800664958765 ],mu2=[
    4.378244151290286 , 74.93711588049547 ] total number of
    iterations= 9
31  mu1=[ 3.049698162970447 , 35.677793710549864 ],mu2=[
    2.777281161152568 , 58.39125094218829 ] total number of
    iterations= 18
```
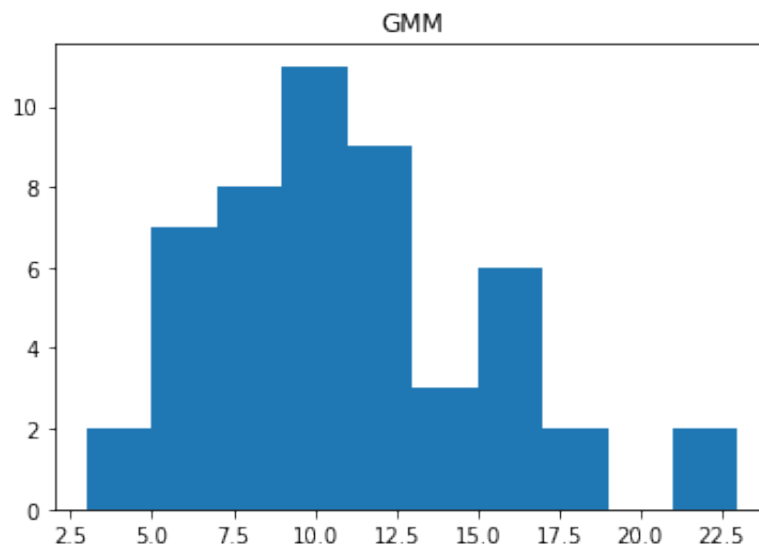
```
32  mu1=[ 0.04830041564780041 , 43.20119999087454 ],mu2=[
    4.820290852579051 , 99.37770639034278 ] total number of
    iterations= 7
33  mu1=[ 2.683818372497481 , 52.69776536228356 ],mu2=[
    5.579858632363873 , 92.0928442357467 ] total number of
    iterations= 7
34  mu1=[ 2.103102347922838 , 39.74264015852395 ],mu2=[
    4.683514312055834 , 72.10140166404754 ] total number of
    iterations= 10
35  mu1=[ 1.9050722120931396 , 45.28556168427721 ],mu2=[
    3.888850058106143 , 95.51868786005926 ] total number of
    iterations= 7
36  mu1=[ 1.217966932603813 , 43.386403151642625 ],mu2=[
    3.9644910623022818 , 69.28408547870853 ] total number of
    iterations= 10
37  mu1=[ 1.4699657079480386 , 45.206730753849 ],mu2=[
    3.4684049139009185 , 64.83729537736575 ] total number of
    iterations= 11
38  mu1=[ 3.772247987679945 , 53.971735638719466 ],mu2=[
    2.6725264456168514 , 61.172305234158266 ] total number of
    iterations= 10
39  mu1=[ 3.484023151232147 , 34.245693344673505 ],mu2=[
    2.2162015261379664 , 56.46513266029492 ] total number of
    iterations= 23
40  mu1=[ 1.6654685051621012 , 55.8290731210131 ],mu2=[
    5.984902395692766 , 78.96466958464332 ] total number of
    iterations= 5
41  mu1=[ 2.116223197482958 , 64.82651145356806 ],mu2=[
    3.329457762612093 , 63.36197863011105 ] total number of
    iterations= 9
42  mu1=[ 3.1358889079370194 , 54.768000392499616 ],mu2=[
    3.849384314442742 , 65.12947860169258 ] total number of
    iterations= 8
43  mu1=[ 0.581119062404257 , 61.63814314559249 ],mu2=[
    3.3661240839425646 , 58.260926370224766 ] total number of
    iterations= 7
44  mu1=[ 1.0382006112690023 , 66.03695087230457 ],mu2=[
    3.208217531208292 , 87.56361837596404 ] total number of
    iterations= 11
45  mu1=[ 3.8022212377216573 , 42.506701341004586 ],mu2=[
    2.760074054725767 , 90.82317636218744 ] total number of
    iterations= 5
```

```
46   mu1=[ 0.39827078755873657 , 62.396269838122656 ],mu2=[
     4.956109112338703 , 93.86805504834047 ] total number of
     iterations= 12
47   mu1=[ 1.135347070198784 , 59.74306249721188 ],mu2=[
     2.871591896482854 , 55.82066991182267 ] total number of
     iterations= 9
48   mu1=[ 1.45475524405298 , 60.66978871723049 ],mu2=[
     2.400002795883992 , 94.22238571061874 ] total number of
     iterations= 12
49   mu1=[ 1.4575503927461755 , 44.062498563433834 ],mu2=[
     3.324706803550839 , 88.27731477686163 ] total number of
     iterations= 5
50   mu1=[ 1.1775053987034951 , 34.99770923313082 ],mu2=[
     3.678004414331871 , 54.29914446593736 ] total number of
     iterations= 21
51   [[ 2.03637021]
52    [54.47833297]]
53   [[ 4.28964582]
54    [79.96791981]]
55   [0.06915319086909583, 33.69625389445514]
56   [0.16998893740998494, 36.049148881677446]
```



GMM

After observing the scatter plot on dataset, randomly choose $\mu_1$, $\mu_2$ in uniform distribution $[U(0, 4), U(30, 70)]$ and $[U(2, 6), U(50, 100)]$, randomly choose $\sigma_1, \sigma_2 \ in \ [[U(1, 20), 0], [0, U(1, 20)]]$ as the initial guess. The distribution of the total number of iterations needed for algorithm to converge is nearly a normal distribution.

## 2.(c) Repeat the task in (b) but with the initial guesses of the parameters generated from the following process:

**i). Run a k-means algorithm over all the data points with K = 2 and label each point with one of the two clusters.**

```python
import numpy as np
import math
import matplotlib.pyplot as plt

#load data
Data_list = []
with open("old.txt", 'r') as in_file:
    for line in in_file.readlines():
        point = []
        point.append(float(line.split()[1]))
        point.append(float(line.split()[2]))
        Data_list.append(point)
Data = np.array(Data_list)
eruptions=[]
waiting=[]
for i in range(len(Data)):
    eruptions.append(Data[i][0])
    waiting.append(Data[i][1])

#k-means and label each point
def distEclud(vecA, vecB):
    return np.sqrt(np.sum(np.power(vecA - vecB, 2)))

def randCent(dataSet, k):
    n = np.shape(dataSet)[1]
    centroids = np.mat(np.zeros([k, n]))
    for j in range(n):
        minj = np.min(dataSet[:,j])
        rangej = float(np.max(dataSet[:,j]) - minj)
        centroids[:,j] = np.mat(minj + rangej * np.random.rand(k, 1))
    print('initial kmeans guess:',centroids)
    return centroids
```

```python
34
35  def KMeans(dataSet, k, distMeans= distEclud, createCent=
    randCent):
36      m = np.shape(dataSet)[0]
37      clusterAssement = np.mat(np.zeros([m,2]))
38      centroids = createCent(dataSet, k)
39      clusterChanged = True
40      while clusterChanged:
41          clusterChanged = False
42          for i in range(m):
43              minDist = np.inf
44              minIndex = -1
45              for j in range(k):
46                  distJ = distMeans(centroids[j,:], dataSet[i,:])
47                  if distJ <  minDist:
48                      minDist = distJ
49                      minIndex = j
50              if clusterAssement[i,0] != minIndex:
51                  clusterChanged = True
52              clusterAssement[i,:] = minIndex, minDist**2
53
54
55          cls0 = dataSet[np.nonzero(clusterAssement[:,0].A == 0)
    [0]]
56          centroids[0,:] = np.mean(cls0, axis = 0)
57          cls1 = dataSet[np.nonzero(clusterAssement[:,0].A == 1)
    [0]]
58          centroids[1,:] = np.mean(cls1, axis = 0)
59      return centroids, cls0,cls1
60
61  center, cls0,cls1 = KMeans(Data, 2)
62  print(center)
```

```
1  initial kmeans guess: [[ 2.69965796 78.79270623]
2   [ 1.91038355 54.62693443]]
3  [[ 4.29793023 80.28488372]
4   [ 2.09433    54.75      ]]
```

The initial guess is $\mu_1$=[ 2.69965796, 78.79270623], $\mu_2$=[ 1.91038355, 54.62693443] and the centre is [ 4.29793023, 80.28488372] and [ 2.09433, 54.75 ].

## ii). Estimate the first guess of the mean and covariance matrices using maximum likelihood over the labeled data points.

```
1  #maximum likelihood over the labeled data points
2  covm0=[]
3  covm1=[]

4  mu0 = cls0.mean(axis=0)

5  sigma0 = (cls0-mu0).T @ (cls0-mu0) / cls0.shape[0]
6  mu1 = cls1.mean(axis=0)

7  sigma1 = (cls1-mu1).T @ (cls1-mu1) / cls1.shape[0]

8  print(mu0)
9  print(mu1)
10 covm0.append(sigma0[0][0])
11 covm0.append(sigma0[1][1])
12 covm1.append(sigma1[0][0])
13 covm1.append(sigma1[1][1])
14 print(covm0)
15 print(covm1)
```

```
1  [ 4.29793023 80.28488372]
2  [ 2.09433 54.75    ]
3  [0.17761716955110854, 31.48279475392103]
4  [0.15427870109999997, 34.4075]
```

mean: [ 4.29793023, 80.28488372],[ 2.09433, 54.75 ] covariance matrix: [[0.17761716955110854, 31.48279475392103][0.15427870109999997, 34.4075]]

## iii). Plot the trajectories of two mean vectors in 2 dimensions (i.e., coordinates vs. iteration).

```python
parameter_dict = {}
parameter_dict["mu1"] = mu0
parameter_dict["sigma1"] = sigma0
parameter_dict["mu2"] = mu1
parameter_dict["sigma2"] = sigma1
parameter_dict["piweight"] = 0.5
parameter_dict["gama_list"] = []


def set_parameter(mu_1, sigma_1, mu_2, sigma_2, pi_weight):
    parameter_dict["mu1"] = mu_1
    parameter_dict["mu1"].shape = (2, 1)
    parameter_dict["sigma1"] = sigma_1
    parameter_dict["mu2"] = mu_2
    parameter_dict["mu2"].shape = (2, 1)
    parameter_dict["sigma2"] = sigma_2
    parameter_dict["piweight"] = pi_weight


def PDF(data, mu, sigma):

    sigma_sqrt = math.sqrt(np.linalg.det(sigma))
    sigma_inv = np.linalg.inv(sigma)
    data.shape = (2, 1)
    mu.shape = (2, 1)
    minus_mu = data - mu
    minus_mu_trans = np.transpose(minus_mu)
    res = (1.0 / (2.0 * math.pi * sigma_sqrt)) * math.exp(
        (-0.5) * (np.dot(np.dot(minus_mu_trans, sigma_inv), minus_mu)))
    return res

def E_step(Data):
    sigma_1 = parameter_dict["sigma1"]
    sigma_2 = parameter_dict["sigma2"]
    pw = parameter_dict["piweight"]
    mu_1 = parameter_dict["mu1"]
    mu_2 = parameter_dict["mu2"]

    parameter_dict["gama_list"] = []
    for point in Data:
        gama_i = (pw * PDF(point, mu_2, sigma_2)) / (
```

```python
                (1.0 - pw) * PDF(point, mu_1, sigma_1) + pw * \
    PDF(point, mu_2, sigma_2))
            parameter_dict["gama_list"].append(gama_i)


def M_step(Data):
    N1 = 0
    N2 = 0
    for i in range(len(parameter_dict["gama_list"])):
        N1 += 1.0 - parameter_dict["gama_list"][i]
        N2 += parameter_dict["gama_list"][i]

    new_mu_1 = np.array([0, 0])
    new_mu_2 = np.array([0, 0])
    for i in range(len(parameter_dict["gama_list"])):
        new_mu_1 = new_mu_1 + Data[i] * (1 - \
    parameter_dict["gama_list"][i]) / N1
        new_mu_2 = new_mu_2 + Data[i] * \
    parameter_dict["gama_list"][i] / N2


    new_mu_1.shape = (2, 1)
    new_mu_2.shape = (2, 1)

    new_sigma_1 = np.array([[0, 0], [0, 0]])
    new_sigma_2 = np.array([[0, 0], [0, 0]])
    for i in range(len(parameter_dict["gama_list"])):
        data_tmp = [0, 0]
        data_tmp[0] = Data[i][0]
        data_tmp[1] = Data[i][1]
        vec_tmp = np.array(data_tmp)
        vec_tmp.shape = (2, 1)
        new_sigma_1 = new_sigma_1 + np.dot((vec_tmp - \
    new_mu_1), (vec_tmp - new_mu_1).transpose()) * (1.0 - \
    parameter_dict["gama_list"][i]) / N1
        new_sigma_2 = new_sigma_2 + np.dot((vec_tmp - \
    new_mu_2), (vec_tmp - new_mu_2).transpose()) * \
    parameter_dict["gama_list"][i] / N2
        # print np.dot((vec_tmp-new_mu_1), (vec_tmp- \
    new_mu_1).transpose())
    new_pi = N2 / len(parameter_dict["gama_list"])
```

```python
        parameter_dict["mu1"] = new_mu_1
        parameter_dict["mu2"] = new_mu_2
        parameter_dict["sigma1"] = new_sigma_1
        parameter_dict["sigma2"] = new_sigma_2
        parameter_dict["piweight"] = new_pi

#Plot the trajectories
def EM_iterate_trajectories(iter_time, Data, mu_1, sigma_1,
 mu_2, sigma_2, pi_weight, esp=0.01):

    mean_trace_1 = [[], []]
    mean_trace_2 = [[], []]

    set_parameter(mu_1, sigma_1, mu_2, sigma_2, pi_weight)
    if iter_time == None:
        while (True):
            old_mu_1 = parameter_dict["mu1"].copy()
            old_mu_2 = parameter_dict["mu2"].copy()
            E_step(Data)
            M_step(Data)
            delta_1 = parameter_dict["mu1"] - old_mu_1
            delta_2 = parameter_dict["mu2"] - old_mu_2

            mean_trace_1[0].append(parameter_dict["mu1"][0]
[0])
            mean_trace_1[1].append(parameter_dict["mu1"][1]
[0])
            mean_trace_2[0].append(parameter_dict["mu2"][0]
[0])
            mean_trace_2[1].append(parameter_dict["mu2"][1]
[0])
            if math.fabs(delta_1[0]) < esp and
math.fabs(delta_1[1]) < esp and math.fabs(
                    delta_2[0]) < esp and
math.fabs(delta_2[1]) < esp:
                break
    else:
        for i in range(iter_time):
            pass

    plt.xlim(xmax=6, xmin=1)
    plt.ylim(ymax=100, ymin=30)
```

```
111        plt.xlabel("eruptions")
112        plt.ylabel("waiting")
113        plt.scatter(eruptions, waiting)
114        plt.plot(mean_trace_1[0], mean_trace_1[1], 'ro')
115        plt.plot(mean_trace_2[0], mean_trace_2[1], 'bo')
116        plt.show()
117
118  def task_1():
119        Mu_1 = np.array([3, 60])
120        Sigma_1 = np.array([[10, 0], [0, 10]])
121        Mu_2 = np.array([1, 30])
122        Sigma_2 = np.array([[10, 0], [0, 10]])
123        Pi_weight = 0.5
124
125        EM_iterate_trajectories(None, Data, Mu_1, Sigma_1, Mu_2,
       Sigma_2, Pi_weight)
126
127  task_1()
```



The choice of input is the same as (b)

**iv). Run your program for 50 times with different initial parameter guesses. Show the distribution of the total number of iterations needed for algorithm to converge.**

```
1  def EM_iterate_times(Data, mu_1, sigma_1, mu_2, sigma_2,
   pi_weight, esp=0.01):
2
```

```python
 3        set_parameter(mu_1, sigma_1, mu_2, sigma_2, pi_weight)
 4        iter_times = 0
 5        while (True):
 6            iter_times += 1
 7            old_mu_1 = parameter_dict["mu1"].copy()
 8            old_mu_2 = parameter_dict["mu2"].copy()
 9            E_step(Data)
10            M_step(Data)
11            delta_1 = parameter_dict["mu1"] - old_mu_1
12            delta_2 = parameter_dict["mu2"] - old_mu_2
13            if math.fabs(delta_1[0]) < esp and
   math.fabs(delta_1[1]) < esp and math.fabs(
14                    delta_2[0]) < esp and math.fabs(delta_2[1]) <
   esp:
15                break
16        return iter_times
17
18  def task_2():
19        center, cls0,cls1 = KMeans(Data, 2)



20        mu0 = cls0.mean(axis=0)



21        sigma0 = (cls0-mu0).T @ (cls0-mu0) / cls0.shape[0]
22        mu1 = cls1.mean(axis=0)



23        sigma1 = (cls1-mu1).T @ (cls1-mu1) / cls1.shape[0]



24        try:
25            Mu_1 = mu0
26            Sigma_1 = sigma0
27            Mu_2 = mu1
28            Sigma_2 = sigma1
29            Pi_weight = 0.5
```

```
30          iter_times = EM_iterate_times(Data, Mu_1, Sigma_1,
    Mu_2, Sigma_2, Pi_weight)
31          print ('input of EM:mu1=[',mu0,'],mu2=[',mu1,']','total
    number of iterations=',iter_times)
32      except Exception:
33          print (Exception)
34      return iter_times
35
36  iteration=[]
37  for i in range(50):
38      iteration.append(task_2())
39  plt.hist(iteration)
40  plt.title('K-means+ML')
41  plt.show
42
```

```
1   initial kmeans guess: [[ 2.68930826 93.75045882]
2    [ 3.65814928 54.44413587]]
3   input of EM:mu1=[ [[ 4.29793023]
4    [80.28488372]] ],mu2=[ [[ 2.09433]
5    [54.75   ]] ] total number of iterations= 4
6   initial kmeans guess: [[ 3.95760765 47.48360481]
7    [ 4.55666981 76.8755731 ]]
8   input of EM:mu1=[ [[ 2.09433]
9    [54.75   ]] ],mu2=[ [[ 4.29793023]
10   [80.28488372]] ] total number of iterations= 4
11  initial kmeans guess: [[ 4.07559629 51.21773633]
12   [ 3.48531779 66.4135674 ]]
13  input of EM:mu1=[ [[ 2.09433]
14   [54.75   ]] ],mu2=[ [[ 4.29793023]
15   [80.28488372]] ] total number of iterations= 4
16  initial kmeans guess: [[ 1.85806105 80.8013877 ]
17   [ 2.34142912 74.71229712]]
18  input of EM:mu1=[ [[ 4.29793023]
19   [80.28488372]] ],mu2=[ [[ 2.09433]
20   [54.75   ]] ] total number of iterations= 4
21  initial kmeans guess: [[ 3.81930558 93.10057294]
22   [ 2.68298853 85.2048207 ]]
23  input of EM:mu1=[ [[ 4.29793023]
24   [80.28488372]] ],mu2=[ [[ 2.09433]
25   [54.75   ]] ] total number of iterations= 4
26  initial kmeans guess: [[ 1.80286387 73.53321619]
```

```
27    [ 3.35986242 74.52127337]]
28  input of EM:mu1=[ [[ 2.09433]
29   [54.75    ]] ],mu2=[ [[ 4.29793023]
30   [80.28488372]] ] total number of iterations= 4
31  initial kmeans guess: [[ 1.71396954 85.04776507]
32   [ 2.47802169 63.06664963]]
33  input of EM:mu1=[ [[ 4.29793023]
34   [80.28488372]] ],mu2=[ [[ 2.09433]
35   [54.75    ]] ] total number of iterations= 4
36  initial kmeans guess: [[ 3.35683763 48.35045179]
37   [ 4.08482155 78.73504159]]
38  input of EM:mu1=[ [[ 2.09433]
39   [54.75    ]] ],mu2=[ [[ 4.29793023]
40   [80.28488372]] ] total number of iterations= 4
41  initial kmeans guess: [[ 2.17308827 61.18730929]
42   [ 1.84641703 74.96256695]]
43  input of EM:mu1=[ [[ 2.09433]
44   [54.75    ]] ],mu2=[ [[ 4.29793023]
45   [80.28488372]] ] total number of iterations= 4
46  initial kmeans guess: [[ 4.97404181 50.19363041]
47   [ 4.29541262 56.94099265]]
48  input of EM:mu1=[ [[ 2.09433]
49   [54.75    ]] ],mu2=[ [[ 4.29793023]
50   [80.28488372]] ] total number of iterations= 4
51  initial kmeans guess: [[ 2.60041207 47.89032733]
52   [ 4.46614127 50.04936223]]
53  input of EM:mu1=[ [[ 2.09433]
54   [54.75    ]] ],mu2=[ [[ 4.29793023]
55   [80.28488372]] ] total number of iterations= 4
56  initial kmeans guess: [[ 1.97204481 90.70224316]
57   [ 1.82401924 49.55252837]]
58  input of EM:mu1=[ [[ 4.29793023]
59   [80.28488372]] ],mu2=[ [[ 2.09433]
60   [54.75    ]] ] total number of iterations= 4
61  initial kmeans guess: [[ 3.46424947 47.69736584]
62   [ 4.96898795 76.70892793]]
63  input of EM:mu1=[ [[ 2.09433]
64   [54.75    ]] ],mu2=[ [[ 4.29793023]
65   [80.28488372]] ] total number of iterations= 4
66  initial kmeans guess: [[ 4.59551223 49.73542565]
67   [ 3.42556647 43.87720278]]
68  input of EM:mu1=[ [[ 4.29793023]
```

```
[80.28488372]] ],mu2=[ [[ 2.09433]
 [54.75     ]] ] total number of iterations= 4
initial kmeans guess: [[ 1.86153908 52.30945672]
 [ 2.1401814  91.44658995]]
input of EM:mu1=[ [[ 2.09433]
 [54.75     ]] ],mu2=[ [[ 4.29793023]
 [80.28488372]] ] total number of iterations= 4
initial kmeans guess: [[ 4.52906079 57.36557084]
 [ 4.51072442 43.14278939]]
input of EM:mu1=[ [[ 4.29793023]
 [80.28488372]] ],mu2=[ [[ 2.09433]
 [54.75     ]] ] total number of iterations= 4
initial kmeans guess: [[ 2.51638293 51.89330513]
 [ 4.10711885 77.54743741]]
input of EM:mu1=[ [[ 2.09433]
 [54.75     ]] ],mu2=[ [[ 4.29793023]
 [80.28488372]] ] total number of iterations= 4
initial kmeans guess: [[ 4.73400128 59.77954756]
 [ 2.86143014 56.08758583]]
input of EM:mu1=[ [[ 4.29793023]
 [80.28488372]] ],mu2=[ [[ 2.09433]
 [54.75     ]] ] total number of iterations= 4
initial kmeans guess: [[ 4.79696047 51.73054456]
 [ 4.47531679 69.03859155]]
input of EM:mu1=[ [[ 2.09433]
 [54.75     ]] ],mu2=[ [[ 4.29793023]
 [80.28488372]] ] total number of iterations= 4
initial kmeans guess: [[ 1.8144436  87.38689893]
 [ 3.21827268 45.39302273]]
input of EM:mu1=[ [[ 4.29793023]
 [80.28488372]] ],mu2=[ [[ 2.09433]
 [54.75     ]] ] total number of iterations= 4
initial kmeans guess: [[ 2.49421492 91.71039159]
 [ 4.83218578 70.43047377]]
input of EM:mu1=[ [[ 4.29793023]
 [80.28488372]] ],mu2=[ [[ 2.09433]
 [54.75     ]] ] total number of iterations= 4
initial kmeans guess: [[ 4.86453894 83.58398761]
 [ 3.44933376 54.0810047 ]]
input of EM:mu1=[ [[ 4.29793023]
 [80.28488372]] ],mu2=[ [[ 2.09433]
 [54.75     ]] ] total number of iterations= 4
```

```
initial kmeans guess: [[ 3.72161422 55.6245174 ]
 [ 4.1645047  94.33519971]]
input of EM:mu1=[ [[ 2.09433]
 [54.75    ]] ],mu2=[ [[ 4.29793023]
 [80.28488372]] ] total number of iterations= 4
initial kmeans guess: [[ 4.72685566 55.22143181]
 [ 3.70937265 57.19244001]]
input of EM:mu1=[ [[ 2.09433]
 [54.75    ]] ],mu2=[ [[ 4.29793023]
 [80.28488372]] ] total number of iterations= 4
initial kmeans guess: [[ 3.96442249 72.11940614]
 [ 4.29090986 91.98123798]]
input of EM:mu1=[ [[ 2.09433]
 [54.75    ]] ],mu2=[ [[ 4.29793023]
 [80.28488372]] ] total number of iterations= 4
initial kmeans guess: [[ 3.41583864 74.9976597 ]
 [ 3.67762104 82.04739125]]
input of EM:mu1=[ [[ 2.09433]
 [54.75    ]] ],mu2=[ [[ 4.29793023]
 [80.28488372]] ] total number of iterations= 4
initial kmeans guess: [[ 2.55211655 50.5522158 ]
 [ 4.09313488 58.99528298]]
input of EM:mu1=[ [[ 2.09433]
 [54.75    ]] ],mu2=[ [[ 4.29793023]
 [80.28488372]] ] total number of iterations= 4
initial kmeans guess: [[ 4.61301796 90.35910438]
 [ 4.25853596 93.48037608]]
input of EM:mu1=[ [[ 2.09433]
 [54.75    ]] ],mu2=[ [[ 4.29793023]
 [80.28488372]] ] total number of iterations= 4
initial kmeans guess: [[ 2.64488169 43.96690331]
 [ 2.93621295 87.03996643]]
input of EM:mu1=[ [[ 2.09433]
 [54.75    ]] ],mu2=[ [[ 4.29793023]
 [80.28488372]] ] total number of iterations= 4
initial kmeans guess: [[ 3.85799134 87.83164891]
 [ 3.82696073 76.04726566]]
input of EM:mu1=[ [[ 4.29793023]
 [80.28488372]] ],mu2=[ [[ 2.09433]
 [54.75    ]] ] total number of iterations= 4
initial kmeans guess: [[ 2.83159657 87.65701463]
 [ 4.21424683 77.71753788]]
```

```
153  input of EM:mu1=[ [[ 4.29793023]
154   [80.28488372]] ],mu2=[ [[ 2.09433]
155   [54.75    ]] ] total number of iterations= 4
156  initial kmeans guess: [[ 1.69719895 95.90672076]
157   [ 3.85157929 66.11071965]]
158  input of EM:mu1=[ [[ 4.29793023]
159   [80.28488372]] ],mu2=[ [[ 2.09433]
160   [54.75    ]] ] total number of iterations= 4
161  initial kmeans guess: [[ 2.21972208 49.80678854]
162   [ 2.5141585  69.92274499]]
163  input of EM:mu1=[ [[ 2.09433]
164   [54.75    ]] ],mu2=[ [[ 4.29793023]
165   [80.28488372]] ] total number of iterations= 4
166  initial kmeans guess: [[ 1.73056185 80.60914335]
167   [ 3.15266925 77.45429361]]
168  input of EM:mu1=[ [[ 4.29793023]
169   [80.28488372]] ],mu2=[ [[ 2.09433]
170   [54.75    ]] ] total number of iterations= 4
171  initial kmeans guess: [[ 3.5227987  71.64158428]
172   [ 4.87262107 63.782147  ]]
173  input of EM:mu1=[ [[ 4.29793023]
174   [80.28488372]] ],mu2=[ [[ 2.09433]
175   [54.75    ]] ] total number of iterations= 4
176  initial kmeans guess: [[ 3.10986969 63.49797944]
177   [ 2.13335105 83.37930532]]
178  input of EM:mu1=[ [[ 2.09433]
179   [54.75    ]] ],mu2=[ [[ 4.29793023]
180   [80.28488372]] ] total number of iterations= 4
181  initial kmeans guess: [[ 4.46790643 45.05398005]
182   [ 2.12275439 54.85599037]]
183  input of EM:mu1=[ [[ 2.09433]
184   [54.75    ]] ],mu2=[ [[ 4.29793023]
185   [80.28488372]] ] total number of iterations= 4
186  initial kmeans guess: [[ 3.63381039 52.25239261]
187   [ 2.14496275 94.66645222]]
188  input of EM:mu1=[ [[ 2.09433]
189   [54.75    ]] ],mu2=[ [[ 4.29793023]
190   [80.28488372]] ] total number of iterations= 4
191  initial kmeans guess: [[ 1.61853871 77.69969726]
192   [ 3.98872359 61.66934164]]
193  input of EM:mu1=[ [[ 4.29793023]
194   [80.28488372]] ],mu2=[ [[ 2.09433]
```

```
 [54.75    ]] ] total number of iterations= 4
initial kmeans guess: [[ 3.71994071 63.39599433]
 [ 3.46833509 84.39146801]]
input of EM:mu1=[ [[ 2.09433]
 [54.75    ]] ],mu2=[ [[ 4.29793023]
 [80.28488372]] ] total number of iterations= 4
initial kmeans guess: [[ 4.16108689 53.33885575]
 [ 2.28946327 45.27875717]]
input of EM:mu1=[ [[ 4.29793023]
 [80.28488372]] ],mu2=[ [[ 2.09433]
 [54.75    ]] ] total number of iterations= 4
initial kmeans guess: [[ 2.90135358 67.89446122]
 [ 2.13383438 74.67634701]]
input of EM:mu1=[ [[ 2.09433]
 [54.75    ]] ],mu2=[ [[ 4.29793023]
 [80.28488372]] ] total number of iterations= 4
initial kmeans guess: [[ 1.90640453 53.61443618]
 [ 3.13891829 77.77733164]]
input of EM:mu1=[ [[ 2.09433]
 [54.75    ]] ],mu2=[ [[ 4.29793023]
 [80.28488372]] ] total number of iterations= 4
initial kmeans guess: [[ 3.43560554 86.47344161]
 [ 2.18061814 79.21997948]]
input of EM:mu1=[ [[ 4.29793023]
 [80.28488372]] ],mu2=[ [[ 2.09433]
 [54.75    ]] ] total number of iterations= 4
initial kmeans guess: [[ 4.1751348  70.45770346]
 [ 4.89781251 92.16961094]]
input of EM:mu1=[ [[ 2.09433]
 [54.75    ]] ],mu2=[ [[ 4.29793023]
 [80.28488372]] ] total number of iterations= 4
initial kmeans guess: [[ 1.81474318 61.91395859]
 [ 4.09298786 56.24930173]]
input of EM:mu1=[ [[ 4.29793023]
 [80.28488372]] ],mu2=[ [[ 2.09433]
 [54.75    ]] ] total number of iterations= 4
initial kmeans guess: [[ 4.6568279  81.10030167]
 [ 4.29498886 46.66900458]]
input of EM:mu1=[ [[ 4.29793023]
 [80.28488372]] ],mu2=[ [[ 2.09433]
 [54.75    ]] ] total number of iterations= 4
initial kmeans guess: [[ 4.27652784 66.5551966 ]
```
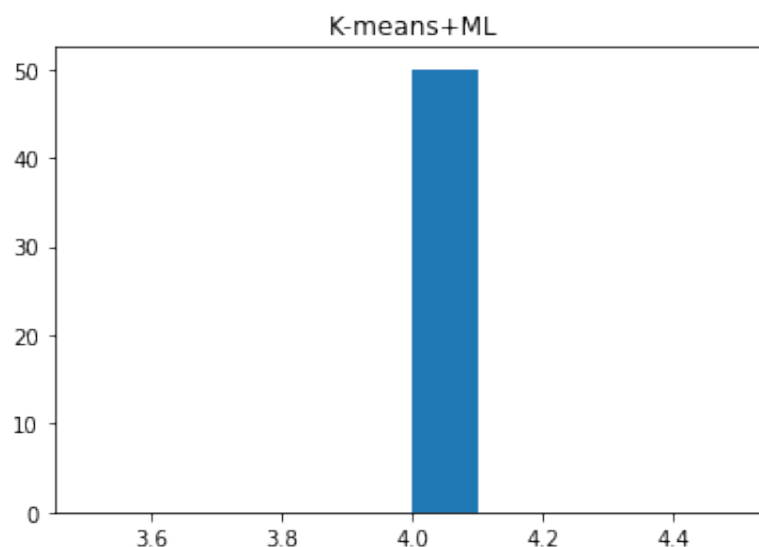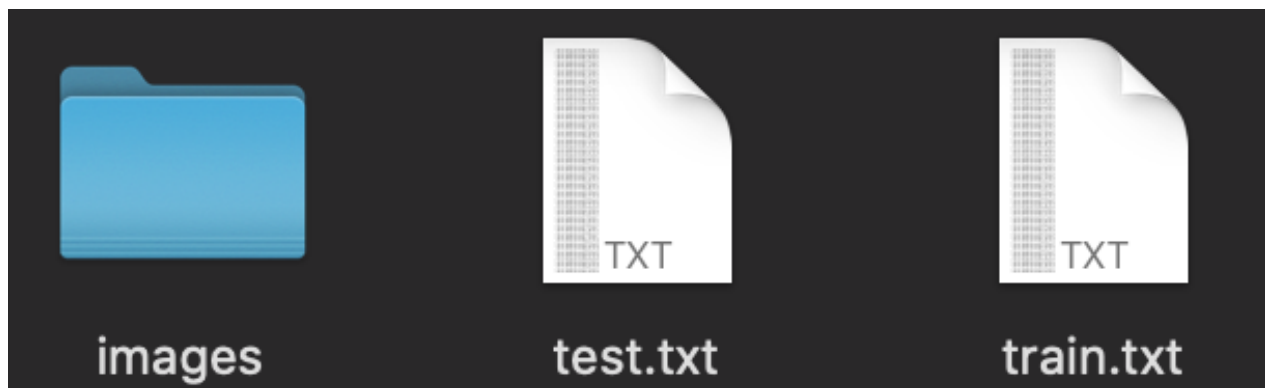
```
237    [ 3.70103244 63.07421508]]
238   input of EM:mu1=[ [[ 4.29793023]
239    [80.28488372]] ],mu2=[ [[ 2.09433]
240    [54.75    ]] ] total number of iterations= 4
241   initial kmeans guess: [[ 2.77564956 70.93614324]
242    [ 5.00486974 81.14260518]]
243   input of EM:mu1=[ [[ 2.09433]
244    [54.75    ]] ],mu2=[ [[ 4.29793023]
245    [80.28488372]] ] total number of iterations= 4
246   initial kmeans guess: [[ 3.36363973 52.05441611]
247    [ 4.57430121 70.03710333]]
248   input of EM:mu1=[ [[ 2.09433]
249    [54.75    ]] ],mu2=[ [[ 4.29793023]
250    [80.28488372]] ] total number of iterations= 4
```

```
1  <function matplotlib.pyplot.show(*args, **kw)>
```



**Compare the algorithm performances of (b) and (c).**

GMM+EM with random initial guess performed depending on its initial guess and the guess' distribution. K-means+GMM+EM performed stable since the centre after k-means calculation is a same result no matter how different the initial input of the k-means. The input of EM is the same, so the iteration times of EM is the same. However, the result of k-means is not perfect, the iteration times is always 4. In (b), with a more randomly initial guess, we have the chance to get a better initial input, so EM could be down within 2 interations.

# Question 3

## 3.(a) Download and unzip The Face Dataset



## 3.(b) Load the training set into a matrix X: there are 540 training images in total, each has 50×50 pixels that need to be concatenated into a 2500-dimensional vector. So the size of X should be 540×2500, where each row is a flattened face image. Pick a face image from X and display that image in grayscale. Do the same thing for the test set. The size of matrix X test for the test set should be 100×2500.

```python
import imageio
import numpy as np
from matplotlib import pylab as plt
import matplotlib.cm as cm
#load data
train_labels, train_data = [], []
for line in open('train.txt'):
    im = imageio.imread(line.strip().split()[0])
    train_data.append(im.reshape(2500,))
    train_labels.append(line.strip().split()[1])
train_data, train_labels = np.array(train_data, dtype=float), np.array(train_labels, dtype=int)
#plot a train image
print(train_data.shape, train_labels.shape)
plt.imshow(train_data[10, :].reshape(50,50), cmap = cm.Greys_r)
plt.show()
```

```
17   test_labels, test_data = [], []
18   for line in open('test.txt'):
19       im = imageio.imread(line.strip().split()[0])
20       test_data.append(im.reshape(2500,))
21       test_labels.append(line.strip().split()[1])
22   test_data, test_labels = np.array(test_data, dtype=float),
     np.array(test_labels, dtype=int)
23   #plot a test image
24   print(test_data.shape, test_labels.shape)
25   plt.imshow(test_data[10, :].reshape(50,50), cmap = cm.Greys_r)
26   plt.show()
```

```
1   (540, 2500) (540,)
```



```
1   (100, 2500) (100,)
```

## 3.(c) Average Face.

```
1  mu = np.mean(train_data, axis= 0)
2  print(mu)
3  plt.imshow(mu.reshape(50,50),cmap = cm.Greys_r)
```
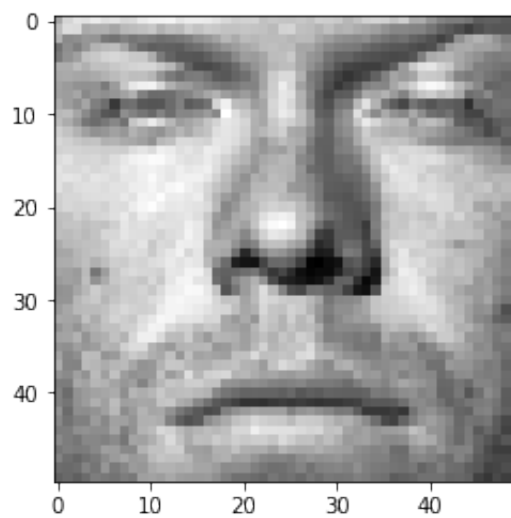
```
1  [59.25185185 56.10185185 52.42222222 ... 67.22222222 64.61851852
2   59.27592593]
```
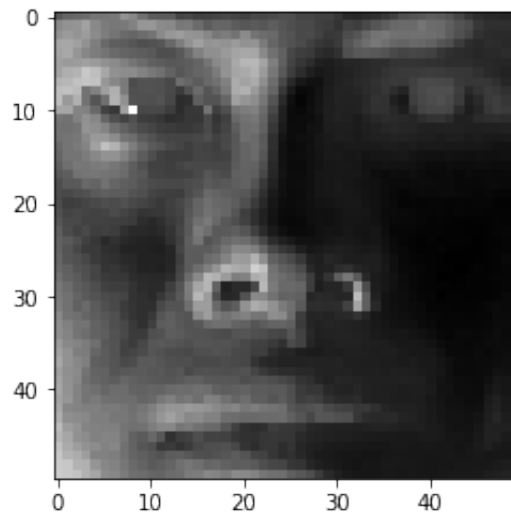
```
1  <matplotlib.image.AxesImage at 0x7faefa248b38>
```

## 3.(d) Mean Subtraction.

```
1  train_sub=train_data-mu
2  test_sub=test_data-mu
3  plt.imshow(train_sub[10, :].reshape(50,50), cmap = cm.Greys_r)
4  plt.show()
5  plt.imshow(test_sub[10, :].reshape(50,50), cmap = cm.Greys_r)
6  plt.show()
```
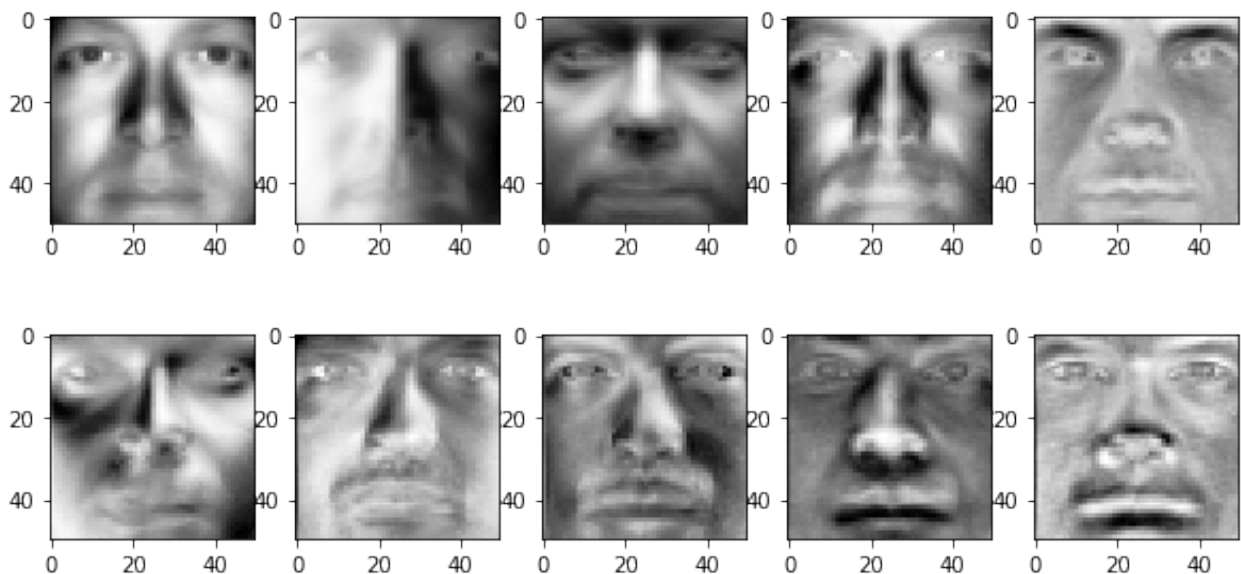
## 3.(e) Eigenface

```
1  X=train_sub
2  u,s,v = np.linalg.svd(X)
3  plt.figure(figsize=(10,5))
4  for i in range(10):
5      plt.subplot(2,5,i+1)
6      plt.imshow(v[i, :].reshape(50,50), cmap = cm.Greys_r)
7  plt.show()
```



## 3.(f) Eigenface Feature.

```
1  X_test=test_sub
2  def Eigenface_Feature(r):
3      V=v[:r,:]
4      F=np.dot(X,V.T)
5
6      V_test=V
7      F_test=np.dot(X_test,V_test.T)
8      return F, F_test
```

## 3.(g) Face Recognition.

```
1   from sklearn.linear_model.logistic import LogisticRegression
2   from sklearn.multiclass import OneVsRestClassifier
3
4   #r=10
5   F,F_test=Eigenface_Feature(10)
6   F_label=train_labels
7   F_test_label=test_labels
8   lr=LogisticRegression(solver='lbfgs',max_iter=10000)
9   ovr=OneVsRestClassifier(lr)
10  ovr.fit(F,F_label)
11  print(ovr.score(F_test,F_test_label))
```

```
1  0.8
```
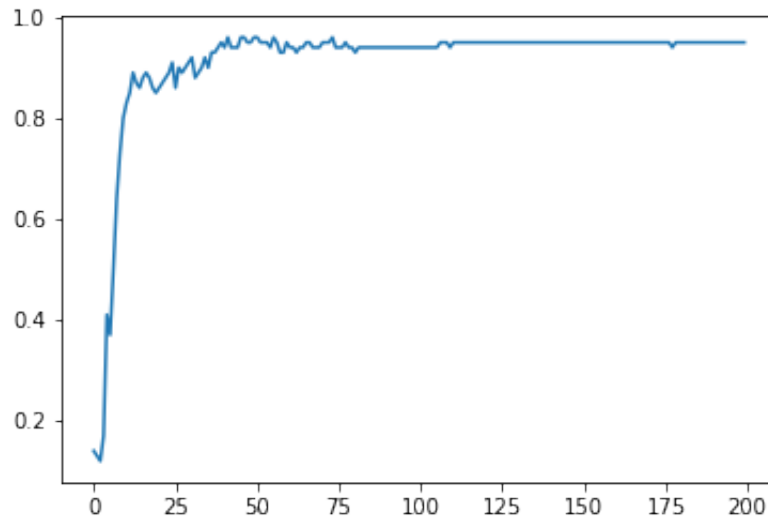
```
1   #r=1-200
2
3   acc=[]
4   R=[]
5   for i in range(200):
6       F,F_test=Eigenface_Feature(i+1)
7       ovr.fit(F,F_label)
8       acc.append(ovr.score(F_test,F_test_label))
9       R.append(i)
10  plt.plot(R,acc)
11  plt.show()
```

# Written Exercises

## Written 1

Question: Show that from the Singular Value Decomposition(SVD) of a matrix $X$, we can obtain the eigendecomposition of $X^T X$. This tells us that we can do an SVD of $X$ and get same result as the eigendecomposition of $X^T X$ but the SVD is faster and easier.

Answer:

The eigenvalue decomposition of a real symmetric matrix can be written as follows:

$$X = Q\Lambda Q^T \tag{1}$$

Assuming that matrix $X$ is an $m * n$ matrix, we can define the SVD of matrix $X$ as follows:

$$X = U\Sigma V^T \tag{2}$$

The columns of $V$ are composed of the eigenvectors of $X^T X$, and the eigenvectors are unit column vectors.

Here's the proof:

$$X = U\Sigma V^T \tag{3}$$

$$\Rightarrow X^T = V\Sigma^T U^T \tag{4}$$

$$\Rightarrow X^T X = V \Sigma^T U^T U \Sigma V = V \Sigma \Sigma^T V^T \tag{5}$$

Note: $\Sigma$ is a $m * n$ matrix, except for the elements on the main diagonal, the value of other elements are all zero. Every element on the main diagonal is a singular value. So we have:

$$\Sigma^T \Sigma = \Sigma^2 \tag{6}$$

Therefore, we can get:

$$\Rightarrow X^T X = V \Sigma^2 V^T \tag{7}$$

Here, $\Sigma$ is a diagonal matrix with singular values as its elements, and $\Sigma \Sigma^T$ is a diagonal matrix consisting of the square of singular values.

Thus, we can see that as a symmetric matrix, the form of the equation (7) of $X^T X$ above is the same as the eigenvalue decomposition of a symmetric matrix.

Therefore, we can know that:

- The nonzero singular values of matrix $X$ is the positive square root of the nonzero eigenvalues of $X^T X$.
- The orthogonal matrix $V$ is the eigenvector of $X^T X$.

In conclusion, we can obtain the eigendecomposition of $X^T X$ from the Singular Value Decomposition(SVD) of a matrix $X$ very fast and easily.

# Written 2

$$z_{il} = x_{il} \Big( \frac{w_l}{\sum_{l=1}^{p} w_l} \Big)^{\frac{1}{2}} \Rightarrow z_{i'l} = x_{i'l} \Big( \frac{w_l}{\sum_{l=1}^{p} w_l} \Big)^{\frac{1}{2}} \tag{8}$$

$$d_e^{(w)}(x_i, x_{i'}) = \sum_{l=1}^{p} (z_{il} - z_{i'l})^2 \tag{9}$$

$$= \sum_{l=1}^{p} [x_{il} \Big( \frac{w_l}{\sum_{l=1}^{p} w_l} \Big)^{\frac{1}{2}} - x_{i'l} \Big( \frac{w_l}{\sum_{l=1}^{p} w_l} \Big)^{\frac{1}{2}}]^2 \tag{10}$$

$$= \sum_{l=1}^{p} [(x_{il} - x_{i'l}) \Big( \frac{w_l}{\sum_{l=1}^{p} w_l} \Big)^{\frac{1}{2}}]^2 \tag{11}$$

$$= \sum_{l=1}^{p} (x_{il} - x_{i'l})^2 \left( \frac{w_l}{\sum_{l=1}^{p} w_l} \right) \tag{12}$$

$$= \frac{\sum_{l=1}^{p} w_l (x_{il} - x_{i'l})^2}{\sum_{l=1}^{p} w_l} \tag{13}$$

# Written 3

```
1   # import modules
2
3   import numpy as np
4   import math
```

```
1   # create the matrix
2
3   M = np.array([[1,0,3],[3,7,2],[2,-2,8],[0,-1,1],[5,8,7]])
4   M = np.mat(M)
5   print (M)
6   print (M.shape)
7   print (type(M))
```

```
1   [[ 1   0   3]
2    [ 3   7   2]
3    [ 2  -2   8]
4    [ 0  -1   1]
5    [ 5   8   7]]
6   (5, 3)
7   <class 'numpy.matrix'>
```

```
1   # (a) Compute the matrices M^T*M and M*M^T
2
3   MT = M.T
4   print (MT)
5   print (MT.shape)
```

```
[[ 1   3   2   0   5]
 [ 0   7  -2  -1   8]
 [ 3   2   8   1   7]]
(3, 5)
```

```
# a-1: Compute the matrix M^T*M

MTM = MT*M
print (MTM)
print (MTM.shape)
```

```
[[ 39   57   60]
 [ 57  118   53]
 [ 60   53  127]]
(3, 3)
```

```
# a-2: Compute the matrix M*M^T

MMT = M*MT
print (MMT)
print (MMT.shape)
```

```
[[ 10    9   26    3   26]
 [  9   62    8   -5   85]
 [ 26    8   72   10   50]
 [  3   -5   10    2   -1]
 [ 26   85   50   -1  138]]
(5, 5)
```

```
1  # (b)&(c)
2
3  # 1.Find the eigenvalues and eigenvectors for M^T*M
4
5  MTM_E, MTM_V = np.linalg.eig(MTM)
6
7  print ("eigenvalues=", MTM_E)
8  print (type(MTM_E))
9  print ("eigenvector=", MTM_V)
10 print (type(MTM_V))
```

```
1  eigenvalues= [2.14670489e+02 9.32587341e-15 6.93295108e+01]
2  <class 'numpy.ndarray'>
3  eigenvector= [[ 0.42615127  0.90453403 -0.01460404]
4   [ 0.61500884 -0.30151134 -0.72859799]
5   [ 0.66344497 -0.30151134  0.68478587]]
6  <class 'numpy.matrix'>
```

```
1  # 1.Find the eigenvalues and eigenvectors for M*M^T
2
3  MMT_E, MMT_V = np.linalg.eig(MMT)
4
5  print ("eigenvalues=", MMT_E)
6  print (type(MMT_E))
7  print ("eigenvector=", MMT_V)
8  print (type(MMT_V))
```

```
1  eigenvalues= [ 2.14670489e+02 -8.88178420e-16  6.93295108e+01
   -3.34838281e-15
2    7.47833227e-16]
3  <class 'numpy.ndarray'>
4  eigenvector= [[-0.16492942 -0.95539856  0.24497323 -0.54001979
   -0.78501713]
5   [-0.47164732 -0.03481209 -0.45330644 -0.62022234  0.30294097]
6   [-0.33647055  0.27076072  0.82943965 -0.12704172  0.2856551 ]
7   [-0.00330585  0.04409532  0.16974659  0.16015949  0.43709105]
8   [-0.79820031  0.10366268 -0.13310656  0.53095405 -0.13902319]]
9  <class 'numpy.matrix'>
```

```
# (d) Method 1

# Find the SVD for the original matrix M from parts (b) and (c)

# d-1: Find the sigma of matrix M
# The result is the same when replacing MMT_E with MTM_E,
  because MTM and MMT have same non-zero eigenvalues, index=0,2
M_sigma = np.array([math.sqrt(MMT_E[0]),math.sqrt(MMT_E[2])])
print ("sigma=", M_sigma)

# check the number of non-zero eigenvalues, = the rank of
  matrix
M_rank = np.linalg.matrix_rank(M)
print (M_rank) # M is rank 2.
```

```
sigma= [14.65163776  8.32643446]
2
```

```
# d-2: Find the VT of matrix M, VT is a 2*2 matrix

rows_V = [0,1] # get row 0 and 1
cols_V = [0,2] # get column 0 and 2
M_V = MTM_V[rows_V,:][:,cols_V]
M_VT = M_V.T
print ("VT=", M_VT)
```

```
VT= [[ 0.42615127  0.61500884]
  [-0.01460404 -0.72859799]]
```

```
# d-3: Find the U of matrix M, U is a 5*2 matrix

cols_U = [0,2] #  get column 0 and 2
M_U = MMT_V[:,cols_U]
print ("U=", M_U)
```

```
1  U= [[-0.16492942  0.24497323]
2   [-0.47164732 -0.45330644]
3   [-0.33647055  0.82943965]
4   [-0.00330585  0.16974659]
5   [-0.79820031 -0.13310656]]
```

```
1  # (e) Method 1
2
3  # e-1: keep only one non-zero singular value, by setting the
   smaller singular value to 0
4
5  sigma_max = max(M_sigma)
6  print (sigma_max)
7  M_sigma_new = np.array([sigma_max])
8  #M_sigma_new = np.mat(M_sigma_new)
9  print (M_sigma_new)
```

```
1  14.651637764976883
2  [14.65163776]
```

```
1   # e-2: Compute the 1D approximation to M
2
3   k=1
4   u,d,vt = M_U[:,:k],M_sigma[:k],M_VT[:,:k][:k,:]
5   print("------U-----")
6   print(u)
7   print("------S-----")
8   print(d)
9   print("------VT-----")
10  print(vt)
11
12  # Compute the 1D approximation to M
13
14  A = np.zeros([1,1])
15  for i in range(1):
16      A[i][i] = d[i]
17  print (A)
```

```
18  tmp = np.dot(u,A)
19  print("1D approximation to M:")
20  print(np.dot(tmp,vt))
```

```
1   ------U-----
2   [[-0.16492942]
3    [-0.47164732]
4    [-0.33647055]
5    [-0.00330585]
6    [-0.79820031]]
7   ------S-----
8   [14.65163776]
9   ------VT-----
10  [[0.42615127]]
11  [[14.65163776]]
12  1D approximation to M:
13  [[-1.02978864]
14   [-2.94487812]
15   [-2.10085952]
16   [-0.02064112]
17   [-4.9838143 ]]
```

```
1   # (d) Method 2
2
3   U,sigma,VT = np.linalg.svd(M)
4
5   print ("U=", U)
6   print ("sigma=", sigma)
7   print ("VT=", VT)
8
9   k=2
10  u,d,vt = U[:,:k],sigma[:k],VT[:,:k][:k,:]
11  print("------U-----")
12  print(u)
13  print("------S-----")
14  print(d)
15  print("------VT-----")
16  print(vt)
```

```
U= [[-0.16492942 -0.24497323  0.9482579   0.09864471
 -0.06214956]
 [-0.47164732  0.45330644 -0.02261948  0.08103373 -0.75165416]
 [-0.33647055 -0.82943965 -0.27341434 -0.18350729 -0.3006445 ]
 [-0.00330585 -0.16974659 -0.14522096  0.97468061  0.00915155]
 [-0.79820031  0.13310656 -0.06671416  0.00505374  0.58368021]]
sigma= [1.46516378e+01 8.32643446e+00 2.99921582e-16]
VT= [[-0.42615127 -0.61500884 -0.66344497]
 [ 0.01460404  0.72859799 -0.68478587]
 [-0.90453403  0.30151134  0.30151134]]
------U-----
[[-0.16492942 -0.24497323]
 [-0.47164732  0.45330644]
 [-0.33647055 -0.82943965]
 [-0.00330585 -0.16974659]
 [-0.79820031  0.13310656]]
------S-----
[14.65163776  8.32643446]
------VT-----
[[-0.42615127 -0.61500884]
 [ 0.01460404  0.72859799]]
```

```python
# (e) Method 2

k=1
u,d,vt = u[:,:k],d[:k],vt[:,:k][:k,:]
print("------U-----")
print(u)
print("------S-----")
print(d)
print("------VT-----")
print(vt)

# Compute the 1D approximation to M
A = np.zeros([1,1])
for i in range(1):
    A[i][i] = d[i]
tmp = np.dot(u,A)
print("1D approximation to M:")
print(np.dot(tmp,vt))
```

```
------U-----
[[-0.16492942]
 [-0.47164732]
 [-0.33647055]
 [-0.00330585]
 [-0.79820031]]
------S-----
[14.65163776]
------VT-----
[[-0.42615127]]
1D approximation to M:
[[1.02978864]
 [2.94487812]
 [2.10085952]
 [0.02064112]
 [4.9838143 ]]
```