

I declare that the assignment submitted on Elearning system is original except for source material explicitly acknowledged, and that the same or related material has not been previously submitted for another course. I also acknowledge that I am aware of University policy and regulations on honesty in academic work, and of the disciplinary guidelines and procedures applicable to breaches of such policy and regulations, as contained in the website <http://www.cuhk.edu.hk/policy/academichonesty/>.

Signed (Student Cao Yuhang) Date: 04/22/2018
Name Cao Yuhang SID 1155092180

Uncompress:

Convert original data to “csv” format. Each row is a 784 dimension data, each column is a pixel feature. And use PCA to reduct 784 dimension to **25** dimension. (I use PCA package of sklearn)

```
1 import os
2 import struct
3 import numpy as np
4 from sklearn.decomposition import PCA
5
6
7 def load_mnist(path, kind='train'):
8     if kind == 'train':
9         labels_path = './ori_data/train-labels-idx1-ubyte'
10        images_path = './ori_data/train-images-idx3-ubyte'
11    elif kind == 'test':
12        images_path = './ori_data/t10k-images-idx3-ubyte'
13        labels_path = './ori_data/t10k-labels-idx1-ubyte'
14
15    with open(labels_path, 'rb') as lbpath:
16        magic, n = struct.unpack('>II', lbpath.read(8))
17        labels = np.fromfile(lbpath, dtype=np.uint8)
18
19    with open(images_path, 'rb') as imgpath:
20        magic, num, rows, cols = struct.unpack('>IIII', imgpath.read(16))
21        images = np.fromfile(imgpath, dtype=np.uint8).reshape(len(labels), 784)
22
23    return images, labels
24
25
26 pca = PCA(n_components=25)
27 X_train, y_train = load_mnist("./")
28 X_train = pca.fit_transform(X_train)
29 X_test, y_test = load_mnist("./", kind='test')
30 X_test = pca.transform(X_test)
31
32 np.savetxt('data/train_img.csv', X_train, fmt='%f', delimiter=',')
33 np.savetxt('data/train_labels.csv', y_train, fmt='%f', delimiter=',')
34 np.savetxt('data/test_img.csv', X_test, fmt='%f', delimiter=',')
35 np.savetxt('data/test_labels.csv', y_test, fmt='%f', delimiter=',')
~
```

Get Initial Centroids:

1. Choose first centroid randomly.
2. Select next initial centroids in such a way that the **Euclidean distance** of that point is **maximum** from other selected initial centroids.
3. Repeat step 2 until we get k initial centroids.
4. The program will ensure initial points contain **at least 7** different labels.

The result is stored in “cluster_id \t centroid \t cluster_number” format, cluster_number represents how many images in this cluster, at beginning, this is just a dump number, just for the consistent in later processing.

```
1 import numpy as np
2
3
4 np.random.seed(0)
5
6 X_train = np.genfromtxt('../data/train_img.csv', dtype=float, delimiter=',')
7 X_label = np.genfromtxt('../data/train_labels.csv', dtype=float, delimiter=',')
8
9 N, D = X_train.shape
10 K = 10
11
12 while True:
13     # initialize
14     initial_centriods = np.zeros((K, D))
15     labels = np.zeros(K)
16     current_indexes = []
17     current_centriods = []
18
19     index = np.random.randint(X_train.shape[0], size=1)[0] # initial point 1
20
21     current_indexes.append(index)
22     current_centriods.append(X_train[index])
23     labels[0] = X_label[index]
24     print ("{}th initial point index: {}, distance is: {}, \
25           label is: {}".format(0, index, 0, X_label[index]))
26
27     for i in range(9):
28         max_distance = -np.inf
29         max_index = None
30         for j in range(X_train.shape[0]):
31             distance = 0
32             for k in range(len(current_centriods)):
33                 distance += np.sum((X_train[j] - current_centriods[k])**2)
34             if distance > max_distance and j not in current_indexes:
35                 max_distance = distance
36                 max_index = j
37             print ("{}th initial point index: {}, distance is: {}, \
38                   label is: {}".format(i+1, max_index, max_distance / len(
39                           current_centriods), X_label[max_index]))
40         current_indexes.append(max_index)
41         current_centriods.append(X_train[max_index])
42         labels[i+1] = X_label[max_index]
43
44     for i in range(10):
45         initial_centriods[i] = current_centriods[i]
46
47     print np.unique(labels)
48     if len(np.unique(labels)) >= 7:
49         break
50
51 with open('initial_point.txt', 'w') as f:
52     for i in range(10):
53         centroid = initial_centriods[i]
54         centroid = ', '.join([str(v) for v in centroid])
55         f.write("{}\t{}\t{}\n".format(i, centroid, 0))
~
```

Q1 (a):

Mapper:

input: split of data

output: "cluster_id \t partial sum of this cluster \t image number of this cluster"

Reducer:

Merge partial sum and image number of same cluster_id, divided sum by number to get mean of this cluster_id

Code and Result:

- mapper1.py: mapper
- reducer1.py: reducer
- format_output.py: get desire output format
- cal_error.py: calculate error between consecutive iterations
- run.sh: run the program in an iterative manner
- ori_centroid_points.txt, old_centroid_points.txt, new_centroid_points.txt: tmp files
- res_a: result of part_a

Code:

Mapper:

```
1 #!/usr/bin/env python
2 import sys
3 import numpy as np
4
5
6 def zero():
7     return 0
8
9
10 def main():
11
12     # read cluster points
13     centroids = np.zeros((10, 25))
14     with open('./old_centroid_points.txt', 'r') as f:
15         i = 0
16         for line in f:
17             index, centroid, counts = line.strip().split("\t")
18             index = int(float(index))
19             centroid = np.fromstring(centroid, sep=',')
20             centroids[index] = centroid
21             i += 1
22
23     local_centroid = np.zeros((10, 25))
24     local_centroid_counts = np.zeros((10, 1))
25
26     for line in sys.stdin:
27         p = np.fromstring(line, sep=',')
28         distances = np.sum((centroids - p) ** 2, axis=1)
29         min_index = np.argmin(distances)
30
31         local_centroid[min_index] += p
32         local_centroid_counts[min_index] += 1
33
34     for i in range(len(local_centroid_counts)):
35         centroid = ','.join([str(v) for v in local_centroid[i]])
36         print str(i) + '\t' + centroid + '\t' + str(local_centroid_counts[i, 0])
37
38
39 if __name__ == '__main__':
40     main()
```

Reducer:

```
1 #!/usr/bin/env python
2
3 import sys
4 import numpy as np
5
6
7 def main():
8     current_centroid_index = None
9     current_counts = 0.
10
11    for line in sys.stdin:
12        centroid_index, local_centroid, local_counts = line.strip().split("\t")
13
14        centroid_index = int(centroid_index)
15        local_counts = float(local_counts)
16        local_centroid = np.fromstring(local_centroid, sep=',')
17
18        if current_centroid_index is None:
19            current_centroid_index = centroid_index
20            current_counts = local_counts
21            current_centroid = local_centroid
22        elif centroid_index != current_centroid_index:
23            if current_counts != 0:
24                current_centroid = current_centroid / current_counts
25                p = ','.join([str(v) for v in current_centroid])
26                print("{}\t{}\t{}".format(current_centroid_index, p, current_counts))
27
28            current_counts = local_counts
29            current_centroid_index = centroid_index
30            current_centroid = local_centroid
31        else:
32            current_centroid += local_centroid
33            current_counts += local_counts
34
35        if current_counts != 0:
36            current_centroid = current_centroid / current_counts
37            p = ','.join([str(v) for v in current_centroid])
38            print("{}\t{}\t{}".format(current_centroid_index, p, current_counts))
39
40 if __name__ == "__main__":
41     main()
```

run.sh:

```
1 hadoop fs -rm -r output
2
3 input='data/train_img.csv'
4 output="output"
5
6 m=4
7 r=2
8
9 error=999999
10 i=0
11 cp ./ori_centroid_points.txt ./old_centroid_points.txt
12 cp ./ori_centroid_points.txt ./new_centroid_points.txt
13 rm error.txt
14
15 while [ $error -gt 1000 ]; do
16 # while [ $i -lt 20 ]; do
17     hadoop jar /usr/hdp/current/hadoop-mapreduce-client/hadoop-streaming.jar \
18         -D mapred.map.tasks=$m \
19         -D mapred.reduce.tasks=$r \
20         -D mapred.compress.map.output=true \
21         -file mapper1.py -mapper mapper1.py \
22         -file reducer1.py -reducer reducer1.py \
23         -file old_centroid_points.txt \
24         -input $input \
25         -output output \
26
27     rm new_centroid_points.txt
28     hadoop fs -get output ./
29     hadoop fs -rm -r output
30     cat output/* >> ./new_centroid_points.txt
31     error="$(python cal_error.py)"
32     echo "$iteration ${i}, error is ${error}" >> error.txt
33     cp new_centroid_points.txt old_centroid_points.txt
34
35     rm -rf output
36     i=$(( i + 1 ))
37 done
```


Q1 (b):

code:

get_label.py: calculate the accuracy

```
1 import numpy as np
2
3
4 centroids = np.zeros((10, 25))
5
6 i = 0
7 with open('./centroid_points.txt', 'r') as f:
8     for line in f:
9         centroid = np.fromstring(line.strip().split("\t")[1], sep=',')
10        centroids[i] = centroid
11
12        i += 1
13
14 predictions = [[] for i in range(10)]
15
16 with open('../data/test_img.csv', 'r') as f:
17    i = 0
18    for line in f:
19        p = np.fromstring(line, sep=',')
20        distances = np.sum((centroids - p) ** 2, axis=1)
21        min_index = np.argmin(distances)
22        min_distances = distances[min_index]
23        predictions[min_index].append((i, min_distances))
24        i += 1
25
26
27 nums = np.zeros((10, 1))
28 for i in range(10):
29     predictions[i] = sorted(predictions[i], key=lambda x: x[1])
30     nums[i] = len(predictions[i])
31
32
33 labels = np.genfromtxt('../data/test_labels.csv', delimiter=',')
34
35 for s in [0.05, 0.1, 0.5, 1]:
36     print("threshold is: {}".format(s))
37
38     correctly_clustered_images = np.zeros(10)
39     accuracies = np.zeros(10)
40
41     res = np.zeros(10000)
42
43     for i in range(10):
44         tmp = np.zeros(10)
45         threshold = s * nums[i]
46         above = False
47
48         for p in predictions[i]:
49             index = p[0]
50             label = int(labels[index])
```

```
51             tmp[label] += 1
52             if tmp[label] >= threshold:
53                 above = True
54
55             # calculate accuracy
56             correct_num = 0
57             for pair in predictions[i]:
58                 j = pair[0]
59                 res[j] = label
60                 if res[j] == labels[j]:
61                     correct_num += 1
62             accuracy = correct_num / float(len(predictions[i]))
63             print("cluster: {}, num: {}, threshold: {}, true label: {}, correctly_clustered_images: {}, accuracy: {}".format(
64                 i, nums[i], threshold, label, nums[i]*accuracy, np.round(accuracy, 3)))
65             correctly_clustered_images[i] = nums[i]*accuracy
66             accuracies[i] = accuracy
67             break
68
69         if not above:
70             label = np.argmax(tmp)
71
72         # calculate accuracy
73         correct_num = 0
74
75         for pair in predictions[i]:
76             j = pair[0]
77             res[j] = label
78             if res[j] == labels[j]:
79                 correct_num += 1
80
81             accuracy = correct_num / float(len(predictions[i]))
82
83             print("cluster id: {}, num: {}, threshold: {}, true label: {}, correctly_clustered_images: {}, accuracy: {}".format(
84                 i, nums[i], threshold, label, nums[i]*accuracy, np.round(accuracy, 3)))
85
86             correctly_clustered_images[i] = nums[i]*accuracy
87             accuracies[i] = accuracy
88
89             print("correctly_clustered_images: {}".format(np.sum(correctly_clustered_images)))
90             print("overall accuracy is: {}".format(np.sum(res==labels) / float(len(labels))))
91             print
```


Best x: 50% and 100%

Explain:

Firstly observe that as x increases, accuracy also increases, if x is small, we rely more on the points that has smaller distance, in some sense we can reduce the noise (points that far away from centroid) influence. However, why accuracy keep decreasing? In fact this is the phenomenon of “**curse of dimensionality**”, as dimension increase, the data points will concentrated more and more in a thin-shell near the surface. Let's look at cluster 3, in cluster 5%, we can find the true label is 9, but the #correctly cluster image is only 376; in cluster 10%, we can find the true label is 7, and the #correctly cluster image is 604, although label 9 is more similar to centroid, in the cluster 3 it is **minor**; although label 7 is far away from centroid, in cluster 3 it is **majority**, this is why accuracy increase.

Q1 (c):

- mapper1.py: mapper
- reducer1.py: reducer
- get_data.py: get desire data folders
- cal_error.py: calculate error between consecutive iterations
- run.sh: run the program in an iterative manner
- ori_centroid_points.txt, old_centroid_points.txt, new_centroid_points.txt: tmp files
- res: result of part 5 folders

Notice: (5_folder and 10_folder are almost same, just modify 5 to 10)

mapper and reducer are same with Q1 (a)

get_data.py: split total data (mixture with training and testing data) into 5 folder, loop over 5 folder, every time choose 1 folder as test set and 4 for train set

```
1 import numpy as np
2 import os
3
4
5 np.random.seed(0)
6
7 X_train = np.genfromtxt('../data/mixture_img.csv', delimiter=',')
8 X_label = np.genfromtxt('../data/mixture_labels.csv', delimiter=',')
9
10 n_folder = 5
11 interval = X_train.shape[0] / n_folder
12
13 # shuffle
14 indices = np.random.permutation(X_train.shape[0])
15 X_train = X_train[indices]
16 X_label = X_label[indices]
17
18 if os.path.isdir('5_folder_input'):
19     os.system("rm -rf 5_folder_input")
20 os.system("mkdir 5_folder_input")
21
22 # split into n folders
23 for i in range(n_folder):
24     # test folder
25     test_folder_num = i
26     test_indices = range(test_folder_num * interval,
27                           (test_folder_num + 1) * interval)
28     test_image = X_train[test_indices]
29     test_label = X_label[test_indices]
30
31     # train folder
32     train_image = np.delete(X_train, test_indices, axis=0)
33     train_label = np.delete(X_label, test_indices)
34
35     if os.path.isdir('5_folder_input/iter_{}'.format(i)):
36         os.system("rm -rf 5_folder_input/iter_{}".format(i))
37     os.mkdir('5_folder_input/iter_{}'.format(i))
38
39     np.savetxt('5_folder_input/iter_{}//train_image.csv'.format(i),
40                train_image, fmt='%.f', delimiter=',')
41     np.savetxt('5_folder_input/iter_{}//train_label.csv'.format(i),
42                train_label, fmt='%.f', delimiter=',')
43     np.savetxt('5_folder_input/iter_{}//test_image.csv'.format(i),
44                test_image, fmt='%.f', delimiter=',')
45     np.savetxt('5_folder_input/iter_{}//test_label.csv'.format(i),
46                test_label, fmt='%.f', delimiter=',')
```

run.sh: run 5 experiments for the corresponding 5 dataset (different test and train folder combination)

```
1 m=4
2 r=2
3
4 n=0
5
6 while [ $n -lt 5 ]; do
7     hadoop fs -rm -r output
8     rm -rf output
9     input="5_folder_input/iter_${n}/train_image.csv"
10    output="output"
11    cp ./ori_centroid_points.txt ./old_centroid_points.txt
12    cp ./ori_centroid_points.txt ./new_centroid_points.txt
13    error=99999
14    rm -rf "iter_${n}_error.txt"
15    i=0
16
17    while [ $error -gt 1000 ]; do
18        hadoop jar /usr/hdp/current/hadoop-mapreduce-client/hadoop-streaming.jar \
19            -D mapred.map.tasks=$m \
20            -D mapred.reduce.tasks=$r \
21            -D mapred.compress.map.output=true \
22            -file mapper1.py -mapper mapper1.py \
23            -file reducer1.py -reducer reducer1.py \
24            -file old_centroid_points.txt \
25            -input $input \
26            -output output \
27
28        rm new_centroid_points.txt
29        hadoop fs -get output ./
30        hadoop fs -rm -r output
31        cat output/* >> ./new_centroid_points.txt
32        error=$(python cal_error.py)
33        echo "$iteration ${}, error is ${error}" >> "iter_${n}_error.txt"
34        cp new_centroid_points.txt old_centroid_points.txt
35
36        rm -rf output
37        i=$((i + 1))
38    done
39    cp new_centroid_points.txt "iter_${n}_centroid_points.txt"
40    n=$((n + 1))
41 done
```

The result is tested on x=100%, which is best x in part (b)

5-folder results:

Testing set	Classification Accuracy
Part 1	0.5931
Part 2	0.5898
Part 3	0.597
Part 4	0.594
Part 5	0.5894
Average	0.59266

10-folder results:

Testing set	Classification Accuracy
Part 1	0.5957
Part 2	0.5893
Part 3	0.5983
Part 4	0.5893
Part 5	0.5947
Part 6	0.6034
Part 7	0.598
Part 8	0.5901
Part 9	0.5849
Part 10	0.5953
Average	0.5939

Comparison between with PCA and without PCA:

Compared with without PCA in homework3, the accuracy is almost the same (all about 0.59~0.6), just **lower** a little. PCA represents each 28x28-pixel image of handwritten digit as the linear combination of 25 principal “eigen-digits”, it **loses some information**, this maybe the reason why accuracy decrease a little, but the performance is greatly improved.

Q1 (b):

Code files

- mapper1.py: mapper1
- reducer1.py: reducer1
- mapper2.py: mapper1
- reducer2.py: reducer1
- get_initial_parameters.py
- run.sh: run the program in an iterative manner
- res: result of part 5 folders

get_initial_parameters.py:

- get initial parameter, I just compute the mean, cov of original data, pi is all initialized to be 0.1

```
1 import numpy as np
2
3
4 K, D = 10, 25
5
6 X = np.genfromtxt('./data/train_img.csv', delimiter=',')
7 N = X.shape[0]
8 interval = N / K
9
10 tmp = np.mean(X, axis=0)
11 mu = np.zeros((K, D))
12 cov = np.zeros((K, D*D))
13
14 for k in range(K):
15     x = X[k * interval: (k+1) * interval]
16     mu[k] = np.mean(x, axis=0).reshape(-1)
17     cov[k] = np.cov(x.T).reshape(-1)
18
19 pi = np.ones(10,)
20 pi = pi / np.sum(pi)
21
22 with open('ori_params.txt', 'w') as f:
23     for k in range(K):
24         m = ','.join([str(v) for v in mu[k]])
25         c = ','.join([str(v) for v in cov[k]])
26         f.write('{}\t{}\t{}\t{}\n'.format(k, m, c, pi[k]))
~
```

First job is to calculate μ_{new}

mapper1.py:

```
#!/usr/bin/env python
import sys
import numpy as np
K, D = 10, 25
def pdf(x, mu, cov):
    D = x.shape[0]
    part1 = 1 / ((2 * np.pi) ** (D / 2) * np.linalg.det(cov) ** 0.5)
    tmp = (x - mu).reshape(-1, 1)
    part2 = -np.matmul(np.matmul(tmp.T, np.linalg.inv(cov)), tmp).item() / 2.0
    return float(part1 * np.exp(part2))
def read_parameters():
    miu = np.zeros((K, D))
    pi = np.zeros((K,))
    cov = np.zeros((K, D*D))
    with open('ori_params.txt', 'r') as f:
        for line in f:
            k, miu_k, cov_k, pi_k = line.strip().split('\t')
            miu_k = np.fromstring(miu_k, sep=',')
            cov_k = np.fromstring(cov_k, sep=',')
            pi_k = float(pi_k)
            k = int(k)
            miu[k] = miu_k
            cov[k] = cov_k
            pi[k] = pi_k
    return miu, cov, pi
def main():
    miu, cov, pi = read_parameters()
    local_miu = np.zeros((K, D)).astype(np.float64)
    local_N = np.zeros((K,)).astype(np.float64)
35,1          Top
```

```
local_N = np.zeros((K,)).astype(np.float64)
for line in sys.stdin:
    x = np.fromstring(line, sep=',')
    numerators = np.zeros((K,))
    # calculate gamma
    for k in range(K):
        numerators[k] = pi[k] * pdf(x, miu[k], cov[k].reshape(D, D))
    gamma = numerators / np.sum(numerators)
    # calculate partial sum
    for k in range(K):
        local_miu[k] += gamma[k] * x
        local_N[k] += gamma[k]
    # emit in form: k local_miu[k] local_N[k]
    for k in range(K):
        m = ','.join([str(v) for v in local_miu[k]])
        print '{}\t{}'.format(k, m, local_N[k])
if __name__ == '__main__':
    main()
~
```

reducer1.py:

```
yuhangcao — CAOYuHang@dic10: ~/asg4/GMM — ssh CAOYuHang@dic10.ie.cuhk.edu.hk — 120x41
1 #!/usr/bin/env python
2
3 import sys
4 import numpy as np
5
6 K, D = 10, 25
7
8
9 def main():
10     current_cluster = None
11     current_miu = 0
12     current_N = 0
13
14     global_miu = np.zeros((K, D))
15     global_N = np.zeros(K)
16
17     for line in sys.stdin:
18         k, tmp = line.strip().split("\t")
19         k = int(k)
20         local_miu, local_N = tmp.split(" ")
21         local_miu = np.fromstring(local_miu, sep=',')
22         local_N = float(local_N)
23
24         if current_cluster is None or current_cluster == k:
25             current_cluster = k
26             current_miu += local_miu
27             current_N += local_N
28         else:
29             if current_N != 0:
30                 miu = current_miu / current_N
31
32                 global_miu[current_cluster] = miu
33                 global_N[current_cluster] = current_N
34
35                 current_cluster = k
36                 current_miu = local_miu
37                 current_N = local_N
38
39         if current_N != 0:
40             miu = current_miu / current_N
41
42             global_miu[current_cluster] = miu
43             global_N[current_cluster] = current_N
44
45             current_miu = 0
46             current_N = 0
47
48             pi = global_N / np.sum(global_N)
49
50             for k in range(K):
51                 miu_k = ", ".join([str(v) for v in global_miu[k]])
52                 pi_k = pi[k]
53                 print "{}\t{}\n\t{}".format(k, miu_k, pi_k)
54
55 if __name__ == "__main__":
56     main()
57
58
59
59
```

1,1 Top

```
yuhangcao — CAOYuHang@dic10: ~/asg4/GMM — ssh CAOYuHang@dic10.ie.cuhk.edu.hk — 120x41
40         miu = current_miu / current_N
41         global_miu[current_cluster] = miu
42         global_N[current_cluster] = current_N
43
44         pi = global_N / np.sum(global_N)
45
46         for k in range(K):
47             miu_k = ", ".join([str(v) for v in global_miu[k]])
48             pi_k = pi[k]
49             print "{}\t{}\n\t{}".format(k, miu_k, pi_k)
50
51
52 if __name__ == "__main__":
53     main()
54
55
56
57
58
59
59
```

46,1 Bot

second job is used to calculate cov and pi

mapper2.py:

```
1 #!/usr/bin/env python
2 import sys
3 import numpy as np
4
5 K, D = 10, 25
6
7
8
9 def pdf(x, mu, cov):
10    D = x.shape[0]
11    part1 = 1 / ((2 * np.pi) ** (D / 2) * np.linalg.det(cov) ** 0.5)
12    tmp = (x - mu).reshape(-1, 1)
13    part2 = -np.matmul(np.matmul(tmp.T, np.linalg.inv(cov)), tmp).item() / 2.0
14    return float(part1 * np.exp(part2))
15
16
17 def read_parameters():
18    miu = np.zeros((K, D))
19    pi = np.zeros(K, )
20    with open('output/part-00000', 'r') as f:
21        for line in f:
22            k, miu_k, pi_k = line.strip().split('\t')
23            k = int(k)
24            miu_k = np.fromstring(miu_k, sep=',')
25            miu[k] = miu_k
26            pi[k] = pi_k
27    return miu, pi
28
29
30 def read_ori_parameters():
31    miu = np.zeros((K, D))
32    pi = np.zeros(K, )
33    cov = np.zeros((K, D*D))
34    with open('ori_params.txt', 'r') as f:
35        for line in f:
36            k, miu_k, cov_k, pi_k = line.strip().split('\t')
37            miu_k = np.fromstring(miu_k, sep=',')
38            cov_k = np.fromstring(cov_k, sep=',')
39            pi_k = float(pi_k)
40
41
42    miu[k] = miu_k
43    cov[k] = cov_k
44    pi[k] = pi_k
45
46
47
48 def main():
49
50    miu, cov, pi = read_ori_parameters()
51    new_miu = read_parameters()[0]
52
53    local_cov = np.zeros((K, D*D)).astype(np.float64)
54    local_N = np.zeros((K, ), dtype=np.float64)
55
56    for line in sys.stdin:
57        x = np.fromstring(line, sep=',')
58        numerators = np.zeros((K, ))
59
60        # calculate gamma
61        for k in range(K):
62            numerators[k] = pdf(x, miu[k], cov[k].reshape(D, D))
63        gamma = numerators / np.sum(numerators)
64
65        # calculate partial sum
66        for k in range(K):
67            tmp = (x - new_miu[k]).reshape(-1, 1)
68            local_cov[k] += gamma[k] * np.matmul(tmp, tmp.T).reshape(-1)
69            local_N[k] += gamma[k]
70
71        # emit in form: k local_N[k]
72        for k in range(K):
73            c = ', '.join([str(v) for v in local_cov[k]])
74            print '{}\t{}\t{}'.format(k, c, local_N[k])
75
76
77 if __name__ == '__main__':
78    main()
```

```
49,0-1          Bot
```

reducer2.py:

```
#!/usr/bin/env python
import sys
import numpy as np
K, D = 10, 25
def read_parameters():
    miu = np.zeros((K, D))
    pi = np.zeros(K,)
    with open('output/part-00000', 'r') as f:
        for line in f:
            k, miu_k, pi_k = line.strip().split('\t')
            k = int(k)
            miu_k = np.fromstring(miu_k, sep=',')
            miu[k] = miu_k
            pi[k] = pi_k
    return miu, pi
def main():
    current_cluster = None
    current_cov = 0
    current_N = 0
    global_miu, pi = read_parameters()
    global_N = np.zeros(K)
    global_cov = np.zeros((K, *D))
    for line in sys.stdin:
        k, tmp = line.strip().split("\t")
        local_cov, local_N = tmp.split(" ")
        k = int(k)
        local_cov = np.fromstring(local_cov, sep=',')
        local_N = float(local_N)
        if current_cluster is None or current_cluster == k:
            current_cluster = k
    "reducer2.py" 68L, 1716C
1,1          Top
```

```
current_cluster = k
current_cov += local_cov
current_N += local_N
else:
    if current_N != 0:
        current_cov = current_cov / current_N
    global_cov[current_cluster] = current_cov
    global_N[current_cluster] = current_N
current_cluster = k
current_cov = local_cov
current_N = local_N
if current_N != 0:
    current_cov = current_cov / current_N
global_cov[current_cluster] = current_cov
global_N[current_cluster] = current_N
for k in range(K):
    miu = ','.join([str(v) for v in global_miu[k]])
    cov = ','.join([str(v) for v in global_cov[k]])
    pi = global_N[k] / np.sum(global_N)
    print "{}\t{}\t{}\t{}".format(k, miu, cov, pi)
if __name__ == "__main__":
    main()
~
~
~
~
~
~
~
```

get_label.py: calculate accuracy, determine label using majority method

```
yuhangcao — CAOYuHang@dic10: ~/asg4/GMM — ssh CAOYuHang@dic10.ie.cuhk.edu.hk — 120x41
1 import numpy as np
2 from scipy.stats import multivariate_normal
3
4
5 K, D = 10, 25
6
7
8 miu = np.zeros((K, D))
9 pi = np.zeros((K,))
10 cov = np.zeros((K, D*D))
11
12 with open('ori_params.txt', 'r') as f:
13     for line in f:
14         k, miu_k, cov_k, pi_k = line.strip().split('\t')
15         miu_k = np.fromstring(miu_k, sep=',')
16         cov_k = np.fromstring(cov_k, sep=',')
17         pi_k = float(pi_k)
18
19         k = int(k)
20         miu[k] = miu_k
21         cov[k] = cov_k
22         pi[k] = pi_k
23
24 models = []
25 for k in range(K):
26     models.append(multivariate_normal(
27         mean=miu[k], cov=cov[k].reshape(D, D)))
28
29 X_test = np.genfromtxt("./data/test_img.csv",
30                        delimiter=",").astype(np.float64)
31 y_test = np.genfromtxt("./data/test_labels.csv",
32                        delimiter=",").astype(np.float64)
33
34
35 predictions = np.zeros(y_test.shape[0])
36 counts = np.zeros((K, K))
37
38 for n in range(X_test.shape[0]):
39     x = X_test[n]
40     tmp = np.zeros((10, )) # likelihood
"get_label.py" 66L, 1724C
1,1          Top
```

```
yuhangcao — CAOYuHang@dic10: ~/asg4/GMM — ssh CAOYuHang@dic10.ie.cuhk.edu.hk — 120x41
40     tmp = np.zeros((10, )) # likelihood
41     for k in range(K):
42         tmp[k] = pi[k] * models[k].pdf(x)
43     cluster = np.argmax(tmp)
44     label = int(y_test[n])
45     counts[cluster, label] += 1
46     predictions[n] = cluster
47
48 cluster = {}
49 for k in range(K):
50     cluster[k] = np.argmax(counts[k])
51     print "cluster:", k
52     print counts[k]
53
54 for n in range(y_test.shape[0]):
55     predictions[n] = cluster[predictions[n]]
56
57 for k in cluster.keys():
58     num = np.sum(counts[k])
59     if num == 0:
60         a = 0
61     else:
62         a = np.max(counts[k]) / num
63     print "cluster: {}, num: {}, threshold is: {}, true label is: {}, correctly_clustered_images: {}, accuracy: {}".format(k, num, num, cluster[k], np.max(counts[k]), a)
64
65 print "accuracy: {}".format(np.sum(predictions == y_test) / float(y_test.shape[0]))
66
~
~
```

Result after training for 20 iterations using GMM:

Table 1. The Accuracy of Clustering Performance with x = 100%					
Cluster Number	# images in the entire cluster	# of images considered (m) when determining the cluster label	Major Label of central images	# correctly clustered images	Classification Accuracy (%)
0	941	941	2	415	0.441
1	1261	1261	4	438	0.347
2	1186	1186	8	734	0.619
3	905	905	0	652	0.720
4	768	768	6	763	0.993
5	986	986	1	984	0.998
6	601	601	5	517	0.860
7	1197	1197	4	464	0.388
8	1160	1160	3	775	0.668
9	995	995	7	504	0.507
Total Set	10000	10000	NA	6246	0.6246

Q2

(a)

Additional parameters:

When predict the desired rating, using the top 2 most similar item

How to handle negative effect of missing ratings:

For a movie m, if missing rating of user u, then fill $r(m, u)$ with the average rating of movie m.

1. Calculate mean of each Movie

	User 1	User 2	User 3	User 4	User 5	Mean
Movie A	1		6	2	4	3.25
Movie B	1	3		1	4	2.25
Movie C	6			4	2	4
Movie D	4	4	2	5		3.75
Movie E	4	5	4	0	3	4
Movie F		4	4	5	1	3.5

2. Normalize each movie by its mean, (**filling missing rating with its average, so after subtract the mean, the missing rating is 0**)

	User 1	User 2	User 3	User 4	User 5
Movie A	-2.25	0	2.75	-1.25	0.75
Movie B	-1.25	0.75	0	-1.25	1.75
Movie C	2	0	0	0	-2
Movie D	0.25	0.25	-1.75	1.25	0
Movie E	0	1	0	0	-1
Movie F	0	0.5	0.5	1.5	-2.5

(c)

Code:

```
yuhangcao — CAOYuHang@dic10: ~/asg4/Q2 — ssh CAOYuHang@dic10.ie.cuhk.edu.hk — 120x41
...— ssh CAOYuHang@dic10.ie.cuhk.edu.hk ...ssh CAOYuHang@dic10.ie.cuhk.edu.hk .. sh CAOYuHang@dic10.ie.cuhk.edu.hk +
```

```
1 import numpy
2
3 numpy.random.seed(0)
4
5
6 def matrix_factorization(R, P, Q, K, alpha=0.0002, beta=0.02):
7     Q = Q.T
8     e = 0
9     last_e = 0
10    step = 0
11    while True:
12        for i in xrange(len(R)):
13            for j in xrange(len(R[i])):
14                if R[i][j] > 0:
15                    eij = R[i][j] - numpy.dot(P[i, :], Q[:, j])
16                    for k in xrange(K):
17                        P[i][k] = P[i][k] + alpha * (
18                            2 * eij * Q[k][j] - beta * P[i][k])
19                        Q[k][j] = Q[k][j] + alpha * (
20                            2 * eij * P[i][k] - beta * Q[k][j])
21
22    e = 0
23    for i in xrange(len(R)):
24        for j in xrange(len(R[i])):
25            if R[i][j] > 0:
26                e = e + pow(R[i][j] - numpy.dot(P[i, :], Q[:, j]), 2)
27
28    print 'step: {}, error: {}'.format(step, e, numpy.abs(e - last_e))
29    if numpy.abs(e - last_e) < 1e-8:
30        break
31    else:
32        last_e = e
33        step += 1
34    # if e < 0.001:
35    #     break
36
37    return P, Q.T
38
39
40 if __name__ == "__main__":
"mf.py" 58L, 1668C
```

1,1 Top

```
yuhangcao — CAOYuHang@dic10: ~/asg4/Q2 — ssh CAOYuHang@dic10.ie.cuhk.edu.hk — 120x41
...— ssh CAOYuHang@dic10.ie.cuhk.edu.hk ...ssh CAOYuHang@dic10.ie.cuhk.edu.hk .. sh CAOYuHang@dic10.ie.cuhk.edu.hk +
```

```
31         break
32     else:
33         last_e = e
34         step += 1
35     # if e < 0.001:
36     #     break
37     return P, Q.T
38
39
40 if __name__ == "__main__":
41     R = numpy.array([[1, 1, 6, 4, 4, 0],
42                     [0, 3, 0, 4, 5, 4],
43                     [6, 0, 0, 2, 4, 4],
44                     [2, 1, 4, 5, 0, 5],
45                     [4, 4, 2, 0, 3, 1]]).astype(float)
46
47     R = numpy.array(R)
48
49     N = len(R)
50     M = len(R[0])
51     K = 2
52
53     P = numpy.random.rand(N, K)
54     Q = numpy.random.rand(M, K)
55
56     nP, nQ = matrix_factorization(R, P, Q, K)
57     nR = numpy.dot(nP, nQ.T)
58     print nR
```

46,0~1 Bot

- new converge condition, error change is smaller than 1e-8
- Prediction rating of User 2 on Movie C: 4.94633897
- Iteration times: 49403

(d)

Mistake: when update Q[k][j], the P[i][k] has been updated, it should use the original P[i][k]

Correct code:

```
1 import numpy
2
3 numpy.random.seed(0)
4
5
6 def matrix_factorization(R, P, Q, K, alpha=0.0002, beta=0.02):
7     Q = Q.T
8     e = 0
9     last_e = 0
10    step = 0
11    while True:
12        for i in xrange(len(R)):
13            for j in xrange(len(R[i])):
14                if R[i][j] > 0:
15                    eij = R[i][j] - numpy.dot(P[i, :], Q[:, j])
16                    for k in xrange(K):
17                        tmp = P[i][k]
18                        t1 = P[i][k] + alpha * (
19                            2 * eij * Q[k][j] - beta * P[i][k])
20                        t2 = Q[k][j] + alpha * (
21                            2 * eij * tmp - beta * Q[k][j])
22                        P[i][k] = t1
23                        Q[k][j] = t2
24                    e = 0
25                    for i in xrange(len(R)):
26                        for j in xrange(len(R[i])):
27                            if R[i][j] > 0:
28                                e = e + pow(R[i][j] - numpy.dot(P[i, :], Q[:, j]), 2)
29                                for k in xrange(K):
30                                    e = e + (beta/2) * (pow(P[i][k], 2) + pow(Q[k][j], 2))
31                    print 'step: {}, error: {}, diff: {}'.format(step, e, numpy.abs(e - last_e))
32                    if numpy.abs(e - last_e) < 1e-8:
33                        break
34                    else:
35                        last_e = e
36                        step += 1
37                    # if e < 0.001:
38                    #     break
39    return P, Q.T
40
41
"v2-mf.py" 60L, 1765C
```

1,1

Top

```
31     print 'step: {}, error: {}, diff: {}'.format(step, e, numpy.abs(e - last_e))
32     if numpy.abs(e - last_e) < 1e-8:
33         break
34     else:
35         last_e = e
36         step += 1
37     # if e < 0.001:
38     #     break
39 return P, Q.T
40
41
42 if __name__ == "__main__":
43     R = numpy.array([[1, 1, 6, 4, 4, 0],
44                     [0, 3, 0, 4, 5, 4],
45                     [6, 0, 0, 2, 4, 4],
46                     [2, 1, 4, 5, 0, 5],
47                     [4, 4, 2, 0, 3, 1]]).astype(float)
48
49 R = numpy.array(R)
50
51 N = len(R)
52 M = len(R[0])
53 K = 2
54
55 P = numpy.random.rand(N, K)
56 Q = numpy.random.rand(M, K)
57
58nP, nQ = matrix_factorization(R, P, Q, K)
59nR = numpy.dot(nP, nQ.T)
60print nR
```

34,1

Bot

Result in next page

- Prediction rating of User 2 on Movie C after modify: 4.9459001
 - Iteration times: 52880
-
- Comparison: The object value is a little lower than the false code, but takes more iterations to converge.
 - Reason I guess: The original code uses the updated $p[i][k]$, this will lead update more for 1 step compare with using original $p[i][k]$, so this will lead fewer iterations.