

Q1)

I use numpy to do SVD in this question

1a)

Code:

```
q1a.py > ...
1  import numpy as np
2
3  # init given data
4  preferences = np.array(
5      [[3, 2, 5, 5, 3, 4],
6       [1, 5, 3, 2, 4, 4],
7       [2, 3, 5, 4, 4, 3],
8       [4, 2, 5, 5, 1, 3],
9       [2, 4, 3, 3, 3, 5],
10      [5, 1, 4, 5, 3, 3],
11      [2, 4, 3, 3, 5, 5]])
12
13 # SVD using numpy library
14 U, s, VT = np.linalg.svd(preferences)
15
16 # make Sigma
17 Sigma = np.diag(s)
18
19 # reduce dimension to a 2-D space
20 U_2d = U[:, :2]
21 Sigma_2d = Sigma[:2, :2]
22 VT_2d = VT[:2, :]
23
24 # Print the U_2d, sigma_2d, and VT_2d matrices in 2-D space
25 print("\nU (2-D):")
26 print(U_2d)
27 print("\nSigma (2-D):")
28 print(Sigma_2d)
29 print("\nV (2-D):")
30 print(VT_2d)
31
```

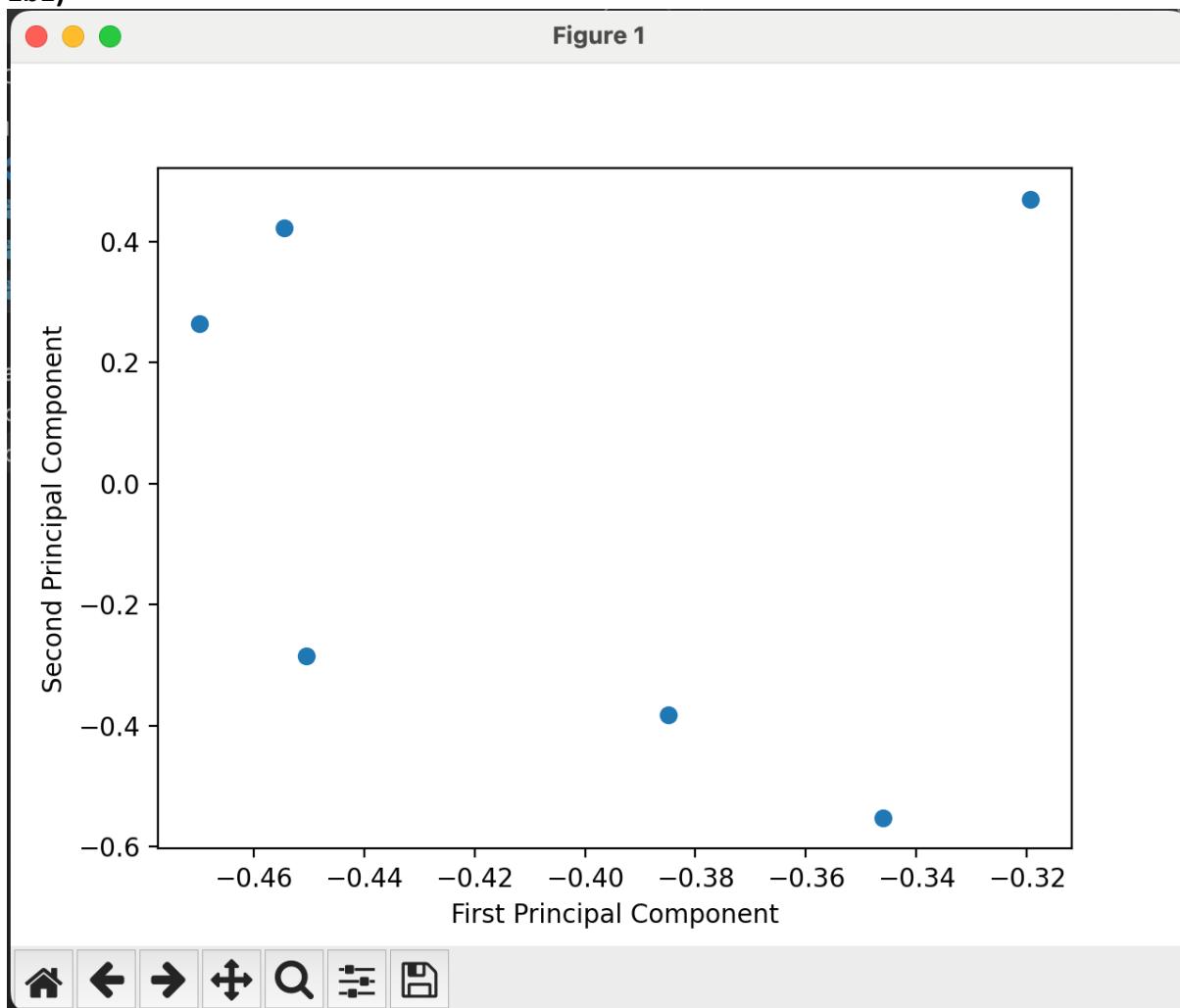
Result:

```
[(base) chansiuchung@Sius-MacBook-Air hw4 % python3 q1a.py
U (2-D):
[[ -0.40783009  0.22897995]
 [ -0.34071273 -0.5221653 ]
 [ -0.38602643 -0.01363104]
 [ -0.36801641  0.46778815]
 [ -0.36251639 -0.27955431]
 [ -0.38009985  0.46689938]
 [ -0.3965387  -0.39969286]]

Sigma (2-D):
[[22.62650088  0.          ]
 [ 0.           6.35359452]]

V (2-D):
[[-0.31940117 -0.3460392 -0.46976188 -0.45444182 -0.38490289 -0.45043972]
 [ 0.46976049 -0.55217132  0.26426701  0.42208251 -0.38165584 -0.28421829]]
```

1b1)



1b2)

Code:

```
q1b2.py > ...
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 # init given data
5 preferences = np.array([
6     [3, 2, 5, 5, 3, 4],
7     [1, 5, 3, 2, 4, 4],
8     [2, 3, 5, 4, 4, 3],
9     [4, 2, 5, 5, 1, 3],
10    [2, 4, 3, 3, 3, 5],
11    [5, 1, 4, 5, 3, 3],
12    [2, 4, 3, 3, 5, 1]])
13
14 # SVD using numpy library
15 U, s, VT = np.linalg.svd(preferences)
16
17 # make Sigma
18 Sigma = np.diag(s)
19
20 # reduce dimension to a 2-D space
21 U_2d = U[:, :2]
22 Sigma_2d = Sigma[:, :2]
23 VT_2d = VT[:, :2, :]
24
25 # new person
26 new_person = np.array([4, 2, 4, 5, 3, 2])
27
28 # representation in the concept space
29 new_person_representation = np.dot(new_person, VT_2d.T)
30
31 print(new_person_representation)
```



Result:

```
[(base) chansiuchung@Sius-MacBook-Air hw4 % python3 q1b2.py
[-8.1765278 2.22877582]]
```

1b3)

Code:

```
q1b3.py > ...
1 import numpy as np
2
3 # init given data
4 preferences = np.array([
5     [3, 2, 5, 5, 3, 4],
6     [1, 5, 3, 2, 4, 4],
7     [2, 3, 5, 4, 4, 3],
8     [4, 2, 5, 5, 1, 3],
9     [2, 4, 3, 3, 3, 5],
10    [5, 1, 4, 5, 3, 3],
11    [2, 4, 3, 3, 5, 1]])
12
13 # SVD using numpy library
14 U, s, VT = np.linalg.svd(preferences)
15
16 # make Sigma
17 Sigma = np.diag(s)
18
19 # reduce dimension to a 2-D space
20 U_2d = U[:, :2]
21 Sigma_2d = Sigma[:, :2]
22 VT_2d = VT[:, :2, :]
23
24 cos_sim = np.dot(preferences[3], preferences[6])/(np.linalg.norm(preferences[3])*np.linalg.norm(preferences[6]))
25 print("cosine similarity:\n{}\n".format(cos_sim))
26
27 cos_sim_2d = np.dot(U_2d[3], U_2d[6])/(np.linalg.norm(U_2d[3])*np.linalg.norm(U_2d[6]))
28 print("\nnew cosine similarity:\n{}\n".format(cos_sim_2d))
```



Result:

```
[(base) chansiuchung@Sius-MacBook-Air hw4 % python3 q1b3.py
cosine similarity:
0.7866066361276135

new cosine similarity:
-0.12246304141824685]]
```

Q2)

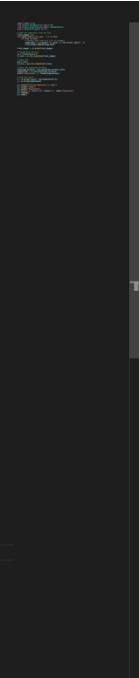
In this question, I use PCA from sklearn, and I followed the suggestion that I should normalize the data points doing PCA so I also used another tool from sklearn to normalize the data before doing the job.

2a1)

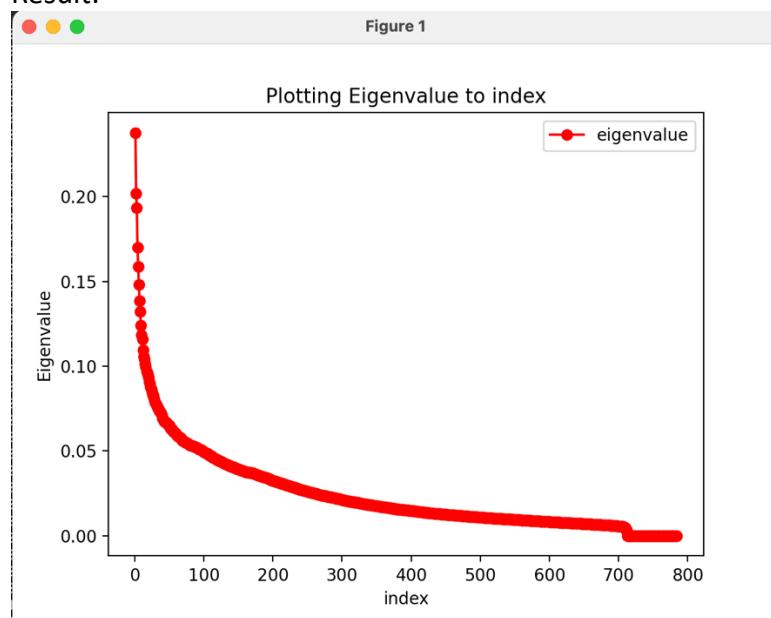
To find the energy left in different number of eigen digits

Code:

```
q2a1.py > ...
1 import numpy as np
2 from sklearn.decomposition import PCA
3 from sklearn.preprocessing import StandardScaler
4 import matplotlib.pyplot as plt
5
6 # Read the image data from the file
7 train_images = []
8 with open("MNIST/train_img", "r") as file:
9     for line in file:
10         # Convert each line to a list of integers
11         image_data = [int(pixel) for pixel in line.strip().split(',')]
12         train_images.append(image_data)
13
14 train_images = np.array(train_images)
15
16 # Standardize the data
17 sc = StandardScaler()
18 X_train = sc.fit_transform(train_images)
19
20 # Apply PCA
21 pca = PCA()
22 X_train = pca.fit_transform(X_train)
23
24 # Check the explained variance
25 explained_variance = pca.explained_variance_ratio_
26 eigenvalues = np.sqrt(explained_variance)
27 print('Eigenvalues: {}'.format(eigenvalues))
28
29 # show the result
30 x = np.array(range(1, len(eigenvalues)+1))
31 y = np.array[eigenvalues]
32
33 plt.title("Plotting Eigenvalue to index")
34 plt.xlabel("index")
35 plt.ylabel("Eigenvalue")
36 plt.plot(x, y, color="red", marker="o", label="eigenvalue")
37 plt.legend()
38 plt.show()
```



Result:



2a2)

First let's find the energy left in different M values. I first normalize the training images for a better result as suggested in doing PCA.

Code:

Change line 21:

```
pca = PCA(n_components=<Eigendigits>)
```

```
q2a1.py > ...
1  import numpy as np
2  from sklearn.decomposition import PCA
3  from sklearn.preprocessing import StandardScaler
4  import matplotlib.pyplot as plt
5
6  # Read the image data from the file
7  train_images = []
8  with open("MNIST/train_img", "r") as file:
9      for line in file:
10          # Convert each line to a list of integers
11          image_data = [int(pixel) for pixel in line.strip().split(',')]
12          train_images.append(image_data)
13
14 train_images = np.array(train_images)
15
16 # Standardize the data
17 sc = StandardScaler()
18 X_train = sc.fit_transform(train_images)
19
20 # Apply PCA
21 pca = PCA(n_components=128)
22 X_train = pca.fit_transform(X_train)
23
24 # Check the explained variance
25 explained_variance = pca.explained_variance_ratio_
26 print("Energy: {}".format(np.sum(explained_variance)))
27 eigenvalues = np.sqrt(explained_variance)
28 print('Eigenvalues: {}'.format(eigenvalues))
29
30 # show the result
31 x = np.array(range(1, len(eigenvalues)+1))
32 y = np.array(eigenvalues)
33
34 plt.title("Plotting Eigenvalue to index")
35 plt.xlabel("index")
36 plt.ylabel("Eigenvalue")
37 plt.plot(x, y, color="red", marker="o", label="eigenvalue")
38 plt.legend()
39 plt.show()
```

Results: (shown later at the end of this question)

To calculate the initial value of the centroid, I use the following code:

```
eigen_to_centroid.py > ...
1  import numpy as np
2  from sklearn.decomposition import PCA
3  from sklearn.preprocessing import StandardScaler
4  import sys
5
6  N = 2
7
8  # Read the image data from stdin
9  train_images = []
10
11 for line in sys.stdin:
12     image_data = [int(pixel) for pixel in line.strip().split(',')]
13     train_images.append(image_data)
14
15 train_images = np.array(train_images)
16
17 # Standardize the data
18 sc = StandardScaler()
19 X_train = sc.fit_transform(train_images)
20
21 # Apply PCA
22 pca = PCA(n_components=N)
23 X_train = pca.fit_transform(X_train)
24
25 # Read centroids
26 centroids = []
27 with open('random_seed_1', 'r') as f:
28     for line in f:
29         centroid = [float(pixel) for pixel in line.strip().split(',')]
30         centroids.append(centroid)
31
32 centroids = np.array(centroids)
33
34 centroid_scaled = sc.transform(centroids)
35
36 centroid_ls = pca.transform(centroid_scaled)
37
38 # Write centroid_ls to a file
39 np.savetxt(f'centroid_ls_{N}.txt', centroid_ls, delimiter=',')
40
41 # Write X_train to a file
42 np.savetxt(f'X_train_{N}.txt', X_train, delimiter=',')
```

Adjust $N = \langle\text{number of eigenvalues}\rangle$ on line 6 to get all 4 pairs of files for Hadoop operation.

And use command:

```
cat MNIST/train_img | python3 eigen_to_centroid.py
```

Now we have 4 initial files for centroid initialization and 4 files of training data after PCA transformation.

Put all 4 $X_{\text{train}}_{\langle\text{num_of_eigenvalues}\rangle}.\text{txt}$ to hdfs and run the following mapreduce code:

1st Mapper:

```
q2a2_mapper.py > ...
1  #!/usr/bin/env python3
2  """q2a2_mapper.py"""
3
4  import sys
5  import numpy as np
6
7  centroid_ls = []
8  vector_ls = []
9
10 with open("centroid_ls_2.txt", 'r') as f:
11     for line in f:
12         #centroid_ls.append(np.array(list(ast.literal_eval(line.strip()))))
13         centroid_ls.append(np.fromstring(line.strip(), sep=","))
14
15 count = 0
16 for line in sys.stdin:
17     count += 1
18     vector_ls.append(np.fromstring(line.strip(), sep=","))

20
21 for vector in vector_ls:
22     # compare
23     nearest_ct_key = -1
24     min_distance = None
25     for key, ct in enumerate(centroid_ls):
26         if min_distance == None:
27             min_distance = np.sum(np.square(vector - ct))
28             nearest_ct_key = key
29         else:
30             cur_dis = np.sum(np.square(vector - ct))
31             if cur_dis < min_distance:
32                 min_distance = cur_dis
33                 nearest_ct_key = key
34
35     print(f"\t{nearest_ct_key}\t{list(vector)})
```

1st reducer:

```
q2a2_reducer.py > ...
1  #!/usr/bin/env python3
2  """q2a2_reducer.py"""
3
4  import sys
5  import numpy as np
6  import ast
7
8  current_centroid_key = None
9  ct_vector = None
10 cur_count = 0
11 centroid_key = None
12
13 for line in sys.stdin:
14     centroid_key, vector = line.strip().split('\t')
15     vector = ast.literal_eval(vector)
16
17     if centroid_key != current_centroid_key:
18         if current_centroid_key != None:
19             print(f"\t{current_centroid_key}\t{cur_count}\t{list(ct_vector/cur_count)}")
20         cur_count = 1
21         current_centroid_key = centroid_key
22         ct_vector = np.array(vector)
23     else:
24         ct_vector += np.array(vector)
25         cur_count += 1
26
27     if current_centroid_key == centroid_key:
28         print(f"\t{current_centroid_key}\t{cur_count}\t{list(ct_vector/cur_count)})
```

2nd Mapper:

```
q2a2_mapper_2.py > ...

7
8 centroid_ls = []
9 vector_ls = []
10
11 count = 0
12 for line in sys.stdin:
13     count += 1
14     vector_ls.append(np.fromstring(line.strip(), sep=""))
15
16 with open("eigen2_output_file", 'r') as f:
17     temp_ls = []
18     for line in f:
19         centroid_id, count ,ct_vector = line.strip().split('\t')
20         temp_ls.append((int(centroid_id), ct_vector))
21
22 temp_ls = sorted(temp_ls, key=lambda x:x[0])
23
24 for i in temp_ls:
25     centroid_ls.append(np.array(ast.literal_eval(i[1])))
26
27 for vector in vector_ls:
28     # compare
29     nearest_ct_key = -1
30     min_distance = None
31     for key, ct in enumerate(centroid_ls):
32         if min_distance == None:
33             min_distance = np.sum(np.square(vector - ct))
34             nearest_ct_key = key
35         else:
36             cur_dis = np.sum(np.square(vector - ct))
37             if cur_dis < min_distance:
38                 min_distance = cur_dis
39                 nearest_ct_key = key
40
41 print(f"{nearest_ct_key}\t{list(vector)}")
```

2nd Reducer (same as 1st reducer):

```
q2a2_reducer.py > ...
1 #!/usr/bin/env python3
2 """q2a2_reducer.py"""
3
4 import sys
5 import numpy as np
6 import ast
7
8 current_centroid_key = None
9 ct_vector = None
10 cur_count = 0
11 centroid_key = None
12
13 for line in sys.stdin:
14     centroid_key, vector = line.strip().split('\t')
15     vector = ast.literal_eval(vector)
16
17     if centroid_key != current_centroid_key:
18         if current_centroid_key != None:
19             print(f"{current_centroid_key}\t{cur_count}\t{list(ct_vector/cur_count)}")
20         cur_count = 1
21         current_centroid_key = centroid_key
22         ct_vector = np.array(vector)
23     else:
24         ct_vector += np.array(vector)
25         cur_count += 1
26
27 if current_centroid_key == centroid_key:
28     print(f"{current_centroid_key}\t{cur_count}\t{list(ct_vector/cur_count)})
```

Loop through 2nd Mapper and reducer until converge using shell command, example:

```
q2a2.sh
1 #!/bin/bash
2
3 for i in {1..150}; do
4     echo "Running iteration $i..."
5     hadoop jar /usr/lib/hadoop-mapreduce/hadoop-streaming-2.10.1.jar \
6     -files hdfs://dicvmc2.ie.cuhk.edu.hk/user/s1155157657/hw4/q2a2_output/eigen2_output_file \
7     -D mapred.map.tasks=5 -D mapred.reduce.tasks=1 -file q2a2_mapper_2.py -mapper q2a2_mapper_2.py \
8     -file q2a2_reducer.py -reducer q2a2_reducer.py -input hw4/X_train_2.txt -output hw4/q2a2_output/eigen2_output_$i
9
10    hadoop fs -cp -f hw4/q2a2_output/eigen2_output_$i/part-00000 hw4/q2a2_output/eigen2_output_file
11    sync
12 done
```

All command used for Hadoop is placed in command_used.sh file

Using the technique of diff the output file, for example:

```
diff <(hadoop fs -cat hw4/q2a2_output/eigen128_output_102/part-00000)
<(hadoop fs -cat hw4/q2a2_output/eigen128_output_103/part-00000)
```

Eigen2 observed converge after 121 iterations

Eigen8 observed converge after 74 iterations.

Eigen32 observed converge after 171 iterations.

Eigen128 observed converge after 102 iterations.

After that I get the output file to local fs using commands:

```
hadoop fs -get hw4/q2a2_output/eigen2_output_121 .
hadoop fs -get hw4/q2a2_output/eigen8_output_75 .
hadoop fs -get hw4/q2a2_output/eigen32_output_172 .
hadoop fs -get hw4/q2a2_output/eigen128_output_103 .
```

For the result, I stored it in the file named: eigen2, eigen8, eigen32, eigen128

For checking accuracy, I use the following program:

Output of the following program will be in this format <prediction>,<actual_label>\t<count>

```
check_details.py > ...
1  #!/usr/bin/env python3
2  """check_details.py"""
3
4  import numpy as np
5  import ast
6
7  vector_ls = []
8  label_ls = []
9  count = 0
10 centroid_ls = []
11
12 with open('X_train_2', 'r') as f:
13     for line in f:
14         count += 1
15         vector_ls.append(np.fromstring(line.strip(), sep=""))
16
17 with open('MNIST/train_label', 'r') as f:
18     for line in f:
19         label_ls.append(int(line.strip()))
20
21 with open('eigen2_output_121', 'r') as f:
22     for line in f:
23         i, each_count, vt = line.strip().split('\t')
24         centroid_ls.append(np.array(ast.literal_eval(vt)))
25
26
27 for i, vector in enumerate(vector_ls):
28     # compare
29     nearest_ct_key = -1
30     min_distance = None
31     for key, ct in enumerate(centroid_ls):
32         if min_distance == None:
33             min_distance = np.sum(np.square(vector - ct))
34             nearest_ct_key = key
35         else:
36             cur_dis = np.sum(np.square(vector - ct))
37             if cur_dis < min_distance:
38                 min_distance = cur_dis
39                 nearest_ct_key = key
40
41     print(f"{str(nearest_ct_key)},{str(label_ls[i])}\t{1}")
```

Then I run the command:

```
python3 check_details.py | sort | python3 check_details_reducer.py
```

outputs:

Eigen2:

```

0, 0    1381
0, 2    44
0, 3    23
0, 4    38
0, 5    29
0, 6    229
0, 7    7
0, 8    36
0, 9    23
1, 0    32
1, 1    6895
1, 2    184
1, 3    589
1, 4    594
1, 5    428
1, 6    657
1, 7    1389
1, 8    723
1, 9    1568
2, 0    899
2, 2    956
2, 3    587
2, 4    445
2, 5    452
2, 6    1866
2, 7    97
2, 8    258
2, 9    148
3, 0    373
3, 1    36
3, 2    1624
3, 3    1698
3, 4    118
3, 5    997
3, 6    1338
3, 7    578
3, 8    436
4, 0    138
4, 1    136
4, 2    116
4, 3    388
4, 4    381
4, 5    389
4, 6    188
4, 7    220
4, 8    287
4, 9    384
5, 0    1811
5, 1    2
5, 2    1382
5, 3    784
5, 4    128
5, 5    597
5, 6    444
5, 7    29
5, 8    562
5, 9    53
6, 0    1872
6, 2    299
6, 3    54
6, 4    4
6, 5    18
6, 6    186
6, 8    27
6, 9    4
7, 0    16
7, 1    4
7, 2    639
7, 3    974
7, 4    864
7, 5    3630
7, 6    71
7, 7    1375
7, 8    473
7, 9    417
8, 0    194
8, 1    451
8, 2    377
8, 3    878
8, 4    565
8, 5    2123
8, 6    1555
8, 7    388
8, 8    1866
8, 9    278
9, 0    15
9, 1    158
9, 2    426
9, 3    3932
9, 4    1234
9, 5    255
9, 6    68
9, 7    2992
9, 8    691
9, 9    1944

```

Eigen8:

```

[1x151517647@elcvmc4 hw4]$ python3 check_details.py | sort | python3 check_details_reducer.py
0, 0    2746
0, 2    14
0, 3    9
0, 4    38
0, 5    43
0, 6    234
0, 7    14
0, 8    44
0, 9    27
1, 0    36
1, 1    6385
1, 2    628
1, 3    11
1, 4    529
1, 5    898
1, 6    494
1, 7    578
1, 8    7977
1, 9    447
2, 0    61
2, 1    13
2, 2    2866
2, 3    323
2, 4    38
2, 5    125
2, 6    854
2, 7    12
2, 8    134
2, 9    15
3, 0    576
3, 1    21
3, 2    338
3, 3    3567
3, 4    1
3, 5    16466
3, 6    48
3, 7    13
3, 8    1417
3, 9    85
4, 0    46
4, 1    5
4, 2    162
4, 3    128
4, 4    171
4, 5    484
4, 6    32
4, 7    488
4, 8    589
4, 9    2592

```

```

5, 0   286
5, 1   382
5, 2   442
5, 3   736
5, 4   358
5, 5   1692
5, 6   71
5, 7   186
5, 8   2545
5, 9   2333
6, 0   477
6, 1   11
6, 2   641
6, 3   154
6, 4   178
6, 5   186
6, 6   4375
6, 7   2
6, 8   37
6, 9   5
7, 0   6
7, 1   1
7, 2   39
7, 3   99
7, 4   497
7, 5   68
7, 6   5
7, 7   1838
7, 8   88
7, 9   483
8, 0   1767
8, 1   10066
8, 2   416
8, 3   59
8, 4   236
8, 5   113
8, 6   5
8, 7   118
8, 8   28
9, 0   2
9, 1   4
9, 2   31
9, 3   69
9, 4   383
9, 5   155
9, 6   3937
9, 7   188
9, 8   2872
9, 9   2844

```

Eigen32:

```

[x1151574570@icvcm4 hw4]$ python3 check_details.py | sort| python3 check_details_reducer.py
0, 0   3342
0, 2   29
0, 3   11
0, 4   42
0, 5   28
0, 6   230
0, 7   18
0, 8   37
0, 9   38
1, 0   31
1, 1   6526
1, 2   699
1, 3   337
1, 4   446
1, 5   421
1, 6   464
1, 7   448
1, 8   1198
1, 9   318
2, 0   71
2, 1   12
2, 2   2881
2, 3   452
2, 4   58
2, 5   168
2, 6   441
2, 7   18
2, 8   143
2, 9   16
3, 0   765
3, 1   24
3, 2   623
3, 3   3613
3, 4   2
3, 5   1814
3, 6   57
3, 7   18
3, 8   2819
3, 9   1541
4, 0   15
4, 1   6
4, 2   158
4, 3   181
4, 4   3382
4, 5   277
4, 6   19
4, 7   1346
4, 8   397
4, 9   2889
5, 0   555
5, 1   144
5, 2   73
5, 3   119
5, 4   777
5, 5   2368
5, 6   111
5, 7   111
5, 8   3876
5, 9   113
6, 0   457
6, 1   13
6, 2   722
6, 3   63
6, 4   141
6, 5   185
6, 6   4579
6, 7   5
6, 8   39
6, 9   6
7, 0   7
7, 1   1
7, 2   34
7, 3   93
7, 4   358
7, 5   63
7, 6   5
7, 7   776
7, 8   78
7, 9   417
8, 0   432
8, 1   7
8, 2   1565
8, 3   974
8, 4   23
8, 5   137
8, 6   11
8, 7   13
8, 8   72
8, 9   7
9, 0   7
9, 1   7
9, 2   35
9, 3   35
9, 4   688
9, 5   58
9, 6   1
9, 7   3527
9, 8   176
9, 9   2844

```

Eigen128:

```
(s1156157@dicvmc4 hw4)$ python3 check_details.py | sort| python3 check_details_reducer.py
0,1      3415
0,2      33
0,3      11
0,4      38
0,5      25
0,6      196
0,7      19
0,8      37
0,9      35
1,0      38
1,1      6527
1,2      696
1,3      616
1,4      455
1,5      440
1,6      399
1,7      449
1,8      1153
1,9      312
2,0      81
2,1      12
2,2      2225
2,3      154
2,4      43
2,5      119
2,6      685
2,7      6
2,8      93
2,9      12
3,0      785
3,1      21
3,2      424
3,3      3733
3,4      4
3,5      1738
3,6      37
3,7      18
3,8      1764
3,9      96
4,0      18
4,1      1
4,2      48
4,3      85
4,4      393
4,5      71
4,6      5
4,7      763
4,8      93
4,9      485
5,0      2559
5,1      146
5,2      124
5,3      198
5,4      329
5,5      2421
5,6      113
5,7      180
5,8      2143
5,9      123
6,0      531
6,1      15
6,2      358
6,3      98
6,4      188
6,5      198
6,6      4421
6,7      2
6,8      33
6,9      5
7,0      3
7,1      3
7,2      24
7,3      78
7,4      452
7,5      56
7,6      1
7,7      3413
7,8      116
7,9      1889
8,0      3880
8,1      11
8,2      1475
8,3      989
8,4      24
8,5      133
8,6      11
8,7      12
8,8      69
8,9      7
9,0      24
9,1      6
9,2      171
9,3      164
9,4      3699
9,5      298
9,6      38
9,7      1491
9,8      386
9,9      3935
```

i) Performance in terms of accuracy:

	Accuracy
Eigen2	15843/60000=26.4%
Eigen8	28369/60000=47.3%
Eigen32	28705/60000=47.8%
Eigen128	29632/60000=49.4%

ii) Energy kept:

Eigendigits	Energy Kept	Graph
2	0.09724988889132913	
8	0.2473398709397076	
32	0.46367542615792245	
128	0.7629398151669453	

iii) In general, the accuracy of the data after PCA is worse than the accuracy before PCA (~63.04%).

The accuracy is better when more principal component vector is used.

2a3)

For a bigger M, the classification accuracy will be higher but the time to run MapReduce will be longer and requires more storage space too.

2b)

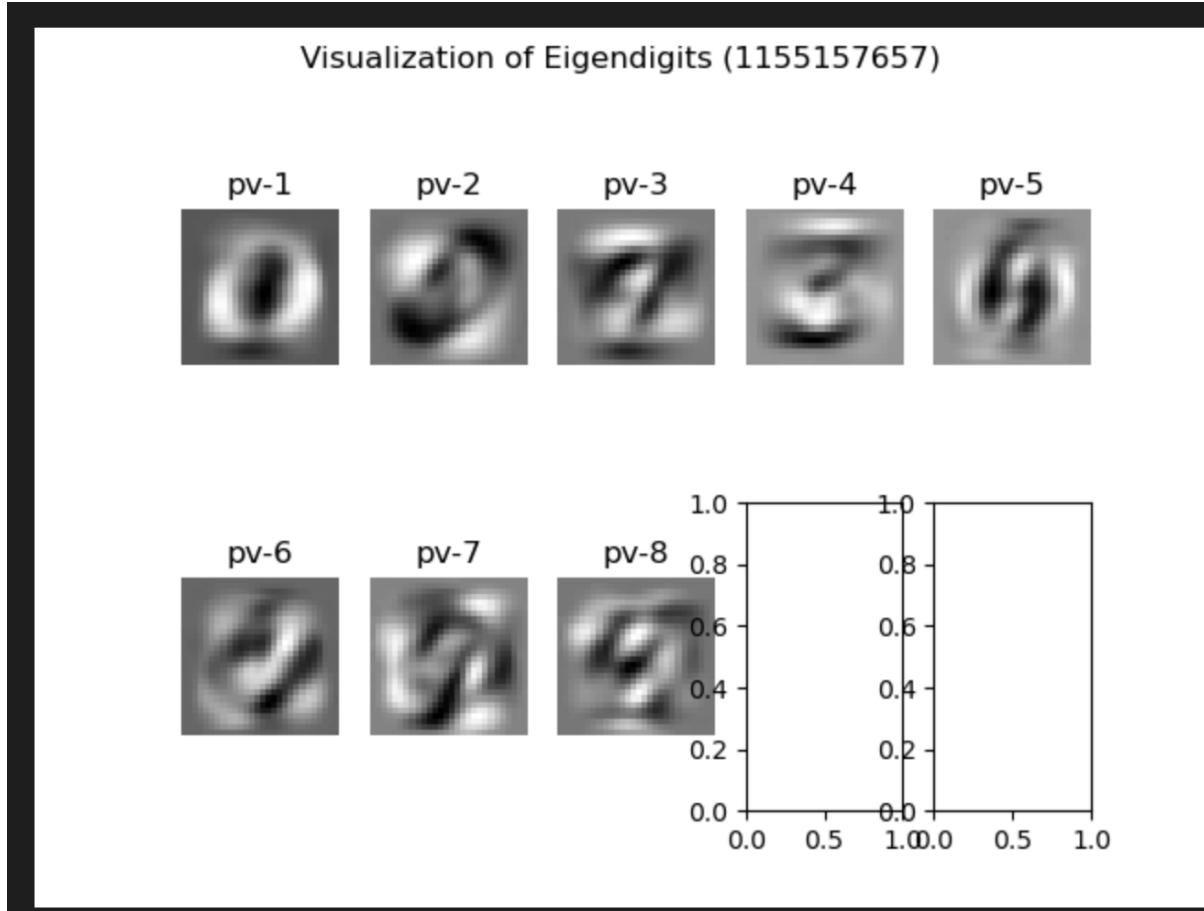
First save the eigen vector to eigen8_vector.txt.

```
find_eigenvectors.py > ...
1 import numpy as np
2 from sklearn.decomposition import PCA
3 from sklearn.preprocessing import StandardScaler
4 import sys
5
6 N = 8
7
8 # Read the image data from stdin
9 train_images = []
10
11 for line in sys.stdin:
12     image_data = [int(pixel) for pixel in line.strip().split(',')]
13     train_images.append(image_data)
14
15 train_images = np.array(train_images)
16
17 # Standardize the data
18 sc = StandardScaler()
19 X_train = sc.fit_transform(train_images)
20
21 # Apply PCA
22 pca = PCA(n_components=N)
23 X_train = pca.fit_transform(X_train)
24
25 # Read centroids
26 centroids = []
27 with open('random_seed_1', 'r') as f:
28     for line in f:
29         centroid = [float(pixel) for pixel in line.strip().split(',')]
30         centroids.append(centroid)
31
32 centroids = np.array(centroids)
33
34 centroid_scaled = sc.transform(centroids)
35
36 centroid_ls = pca.transform(centroid_scaled)
37
38 np.savetxt(f'eigen8_vector.txt', pca.components_, delimiter=',')
```

Second, use the provide code, modify the SID and the pv_dir and fix reduce_dim to 8.

```
viz_principal_vectors.py > ...
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 reduce_dim = 8
5 num_viz_pca = 10
6
7 def main():
8     # load principle vectors from your local machine, e.g., here we suppose that the PCA vectors of the training
9     # set is stored in a file named 'pca_components_20' in the ./mnist/ directory
10    pv_dir = './eigen8_vector.txt'
11    pvs = np.empty((reduce_dim, 784))
12    with open(pv_dir, 'r') as f:
13        # Some pre-processing steps, depending on how you store your principle vectors
14        # Here, we suppose that each row is corresponding to a principle vectors with string
15        # type, e.g., a row can be: '0.5, 0.2, 0.3, ..., 0.21', each element is separated by ','.
16        # for k, line in enumerate(f.readlines()):
17        #     line = line.strip(',')
18        #     # the k-th principle vectors
19        #     pv_k = np.array(list(map(float, line)))
20        #     pvs [k] = pv_k
21
22
23    # visualize the principle vectors
24    fig, ax = plt.subplots(num_viz_pca // 5, 5)
25    # YOU NEED TO FILL IN YOUR SID IN THE FOLLOWING LINE
26    fig.suptitle('Visualization of Eigendigits ({})'.format('1155157657'))
27    ax = ax.reshape(-1)
28    for k, pv in enumerate(pvs [:num_viz_pca]):
29        pv = pv.reshape(28, 28)
30        ax[k].imshow(pv, cmap='gray')
31        ax[k].set_title('pv-{}'.format(k + 1))
32        ax[k].axis('off')
33
34    fig.savefig('./viz_pca_vectors_{}.png'.format(num_viz_pca))
35
36
37 if __name__ == '__main__':
38     main()
```

Output:



Observation:

Some of the eigenvectors looks like numbers , for example '3' and '7' and '9'. Some of them contains common components of a number such as straight lines and circles.

Q3)

I am using the method stated in slide p.38

3a) Code (User-to-User):

```
q3a_user.py > ...
1 import numpy as np
2 from numpy.linalg import norm
3
4 rating_mat = np.array([
5     [5, 0, 4, 1, 4, 2],
6     [2, 3, 4, 5, 0, 0],
7     [0, 2, 0, 3, 5, 4], # Line 7
8     [2, 1, 5, 4, 3, 0],
9     [0, 4, 1, 0, 3, 2]
10])
11
12 # User-User
13 avg_rating_mat = []
14
15 for i in rating_mat:
16     mean_ = i[np.nonzero(i)].mean()
17     avg_rating_mat.append(np.where(i!=0, i-mean_, i))
18
19 avg_rating_mat = np.array(avg_rating_mat)
20
21 cosine_sim = []
22 for key, i in enumerate(avg_rating_mat):
23     cosine = np.dot(avg_rating_mat[2], i)/(norm(avg_rating_mat[2])*norm(i))
24     cosine_sim.append((cosine, key))
25
26 cosine_sim.sort(reverse=True, key=lambda x: x[0])
27
28 # cosine_sim is a list that is reversely sorted with
29 # cosine similarity of each vector to the target vector (cosine_similarity, key_of_vector)
30 print("User-to-User:{(cosine_sim[1][0]*rating_mat[cosine_sim[1][1]][0] + cosine_sim[2][0]*rating_mat[cosine_sim[2][1]][0])/(cosine_sim[1][0]+cosine_sim[2][0])}")
31
```

Code (item-to-item):

```
q3a_item.py > ...
1 import numpy as np
2 from numpy.linalg import norm
3
4 # item-item
5 rating_mat = np.array([
6     [5, 0, 4, 1, 4, 2],
7     [2, 3, 4, 5, 0, 0],
8     [0, 2, 0, 3, 5, 4],
9     [2, 1, 5, 4, 3, 0],
10    [0, 4, 1, 0, 3, 2]
11])
12
13 rating_mat = rating_mat.T
14
15 # Item-Item
16 avg_rating_mat = []
17
18 for i in rating_mat:
19     mean_ = i[np.nonzero(i)].mean()
20     avg_rating_mat.append(np.where(i!=0, i-mean_, i))
21
22 avg_rating_mat = np.array(avg_rating_mat)
23
24 cosine_sim = []
25 for key, i in enumerate(avg_rating_mat):
26     cosine = np.dot(avg_rating_mat[0], i)/(norm(avg_rating_mat[0])*norm(i))
27     cosine_sim.append((cosine, key))
28
29 cosine_sim.sort(reverse=True, key=lambda x: x[0])
30
31 # cosine_sim is a list that is reversely sorted with
32 # cosine similarity of each vector to the target vector (cosine_similarity, key_of_vector)
33 print("Item-to-Item:{(cosine_sim[1][0]*rating_mat[cosine_sim[1][1]][2]+cosine_sim[2][0]*rating_mat[cosine_sim[2][1]][2])/(cosine_sim[1][0]+cosine_sim[2][0])}")
34
```

3ai and 3aii) output:

```
[(base) chansiuchung@Sius-MacBook-Air hw4 % python3 q3a_user.py
User-to-User:3.1865773776947957
[(base) chansiuchung@Sius-MacBook-Air hw4 % python3 q3a_item.py
Item-to-Item:3.8828894492238284]
```

3b)

Code:

```
q3b.py > ...
1 import numpy
2
3 def matrix_factorization(R, P, Q, K, steps=5000, alpha=0.0002, beta=0.02):
4     Q = Q.T
5     for step in range(steps):
6         for i in range(len(R)):
7             for j in range(len(R[i])):
8                 if R[i][j] > 0:
9                     eij = R[i][j] - numpy.dot(P[i,:],Q[:,j])
10                    for k in range(K):
11                        P[i][k] = P[i][k] + alpha * (2 * eij * Q[k][j] - beta * P[i][k])
12                        Q[k][j] = Q[k][j] + alpha * (2 * eij * P[i][k] - beta * Q[k][j])
13    eR = numpy.dot(P,Q)
14    e = 0
15    for i in range(len(R)):
16        for j in range(len(R[i])):
17            if R[i][j] > 0:
18                e = e + pow(R[i][j] - numpy.dot(P[i,:],Q[:,j]), 2)
19                for k in range(K):
20                    e = e + (beta/2) * (pow(P[i][k],2) + pow(Q[k][j],2))
21    if e < 0.001:
22        break
23    return P, Q.T
24
25 R = [
26     [5, 0, 4, 1, 4, 2],
27     [2, 3, 4, 5, 0, 0],
28     [0, 2, 0, 3, 5, 4],
29     [2, 1, 5, 4, 3, 0],
30     [0, 4, 1, 0, 3, 2]
31 ]
32
33 R = numpy.array(R)
34
35 N = len(R)
36 M = len(R[0])
37 K = 2
38
39 P = numpy.random.rand(N,K)
40 Q = numpy.random.rand(M,K)
41
42 nP, nQ = matrix_factorization(R, P, Q, K)
43 nR = numpy.dot(nP, nQ.T)
44
45 print(nR)
```

Result:

```
[(base) chansiuchung@Sius-MacBook-Air hw4 % python3 q3b.py
[[ 4.86056227  4.48916395  3.23429429  0.97645988  4.28725968  2.68425545]
 [ 2.43476992  2.09572284  4.81381808  4.24018163  4.18338121  3.97955249]
 [ 3.02113264  2.66943339  4.53302056  3.56990306  4.27287898  3.74978736]
 [ 1.59467407  1.31380329  4.38056204  4.21910978  3.5225417  3.61936348]
 [ 4.1779845   3.91404453  1.62570427 -0.51652239  2.94933102  1.35487044]]
```

Prediction of User 3 on product A is 3.02113264 while in part a the result is 3.1865776947957 for user-to-user collaborative filtering and 3.8828894492238284 for item-to-item collaborative filtering. The score is lower than the method of collaborative filtering.

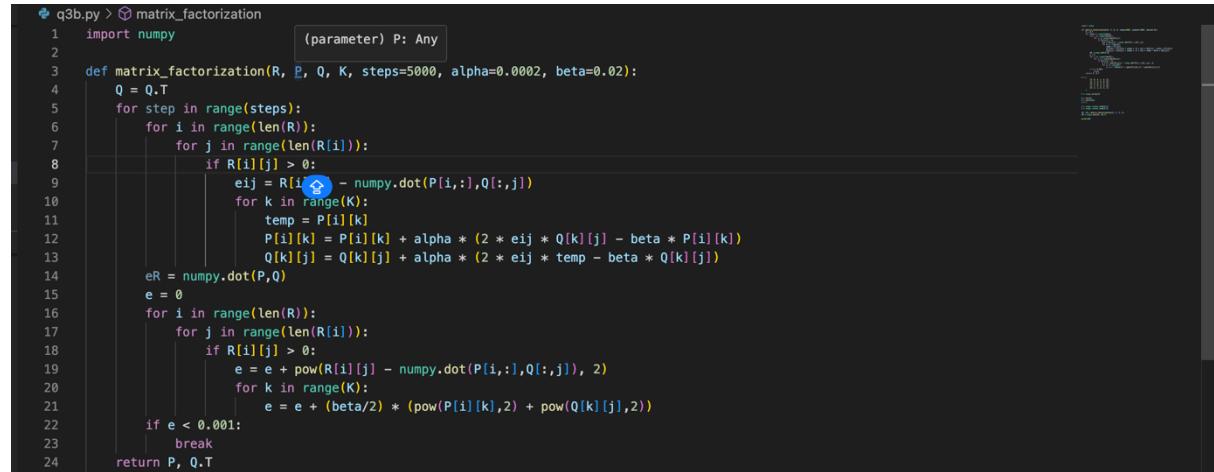
3c)

Code:

$P[i][k]$ is updated first and used in $Q[k][j]$ to calculate new value which is a bug because it affected the accuracy.

Corrected it using a temp variable to store the value first and use it in line 13

Code:



```
q3b.py > matrix_factorization
1 import numpy
2
3 def matrix_factorization(R, P, Q, K, steps=5000, alpha=0.0002, beta=0.02):
4     Q = Q.T
5     for step in range(steps):
6         for i in range(len(R)):
7             for j in range(len(R[i])):
8                 if R[i][j] > 0:
9                     eij = R[i][j] - numpy.dot(P[i,:],Q[:,j])
10                    for k in range(K):
11                        temp = P[i][k]
12                        P[i][k] = P[i][k] + alpha * (2 * eij * Q[k][j] - beta * P[i][k])
13                        Q[k][j] = Q[k][j] + alpha * (2 * eij * temp - beta * Q[k][j])
14
15    eR = numpy.dot(P,Q)
16    e = 0
17    for i in range(len(R)):
18        for j in range(len(R[i])):
19            if R[i][j] > 0:
20                e = e + pow(R[i][j] - numpy.dot(P[i,:],Q[:,j]), 2)
21                for k in range(K):
22                    e = e + (beta/2) * (pow(P[i][k],2) + pow(Q[k][j],2))
23
24    if e < 0.001:
25        break
26
27 return P, Q.T
```

Result:

```
[(base) chansiuchung@Sius-MacBook-Air hw4 % python3 q3b.py
[[4.89300497 5.11716277 3.20631652 0.99268673 4.35212632 2.53231529]
 [2.43987868 2.211466 4.77928555 4.26920415 4.31464542 4.38690432]
 [2.7636278 2.62020556 4.33550825 3.55651084 4.1603402 3.91014088]
 [1.56194123 1.28201502 4.30957074 4.21638527 3.60494843 4.03625474]
 [3.65252452 3.88146515 1.81743177 0.05747216 2.8603837 1.32449939]]
```

The updated value is 2.7636278 which is lower than that in 3b