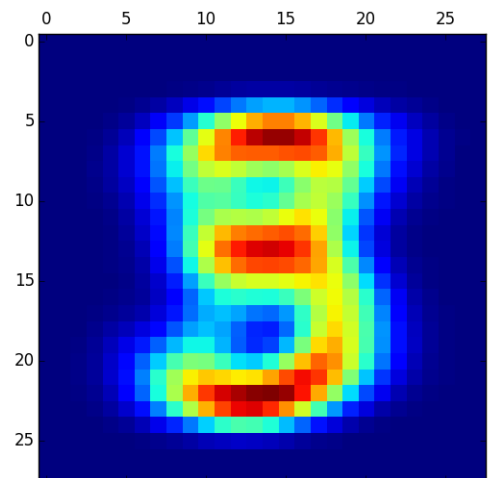
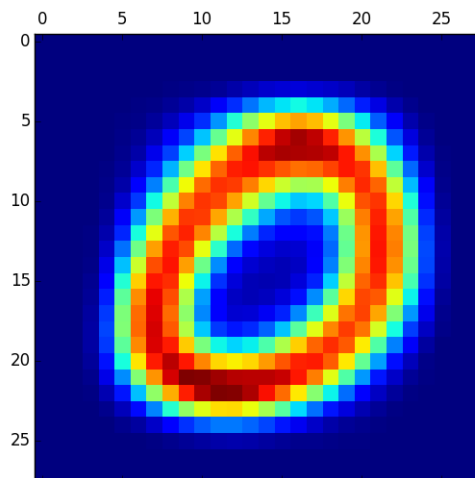


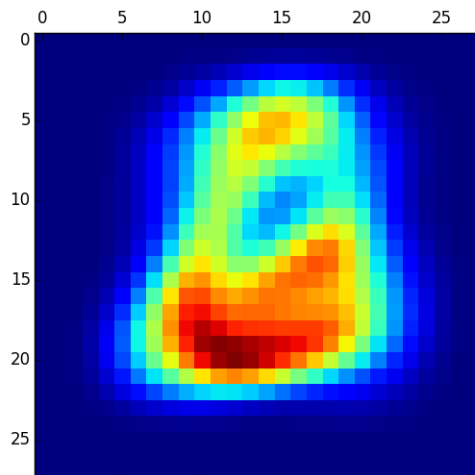
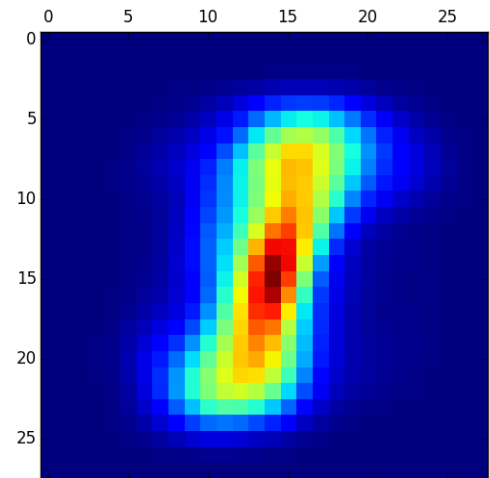
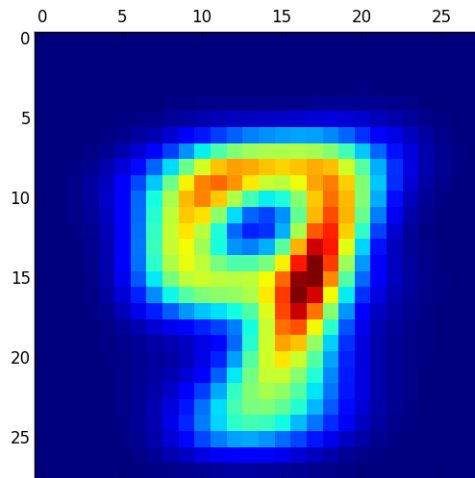
CS189–FALL 2015 — Homework 7 Write up

ZUBO GU, SID 25500921, gu.zubo@berkeley.edu

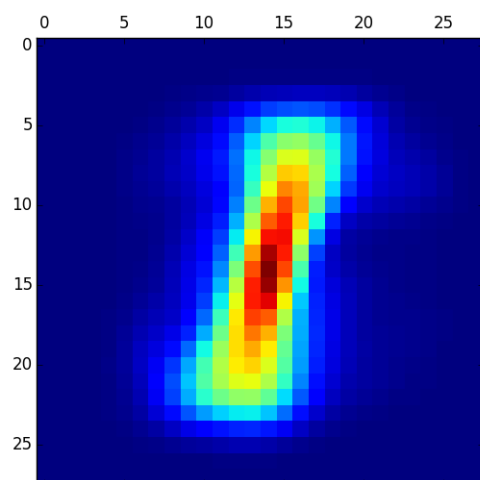
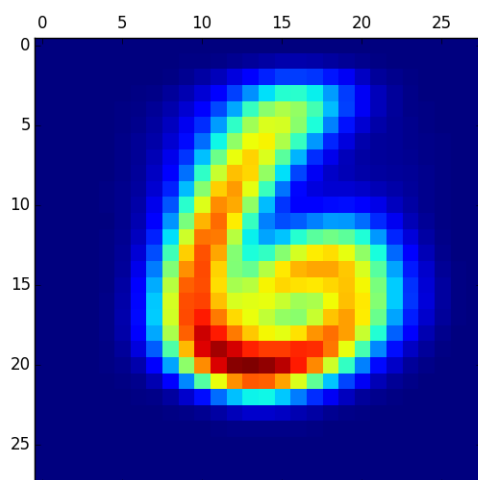
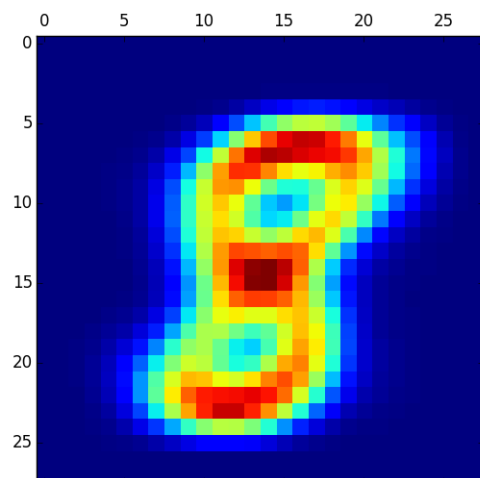
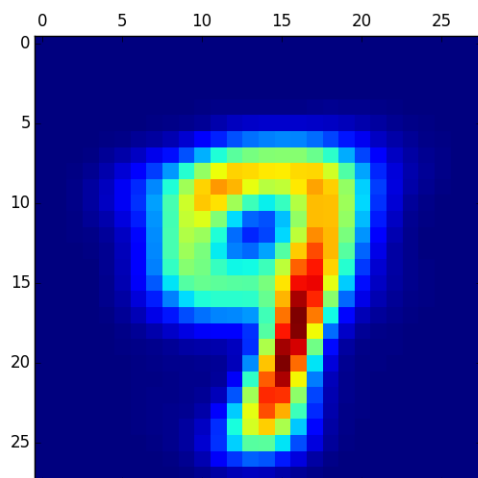
Problem 1

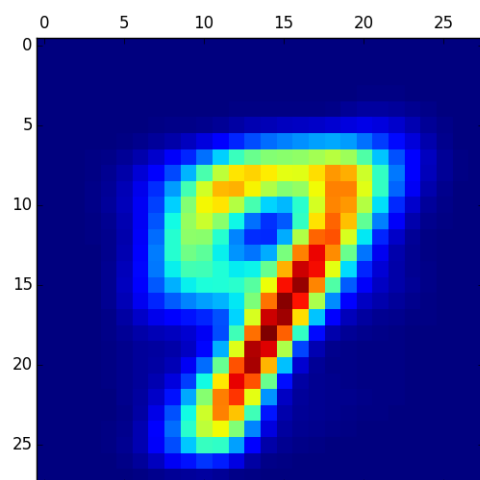
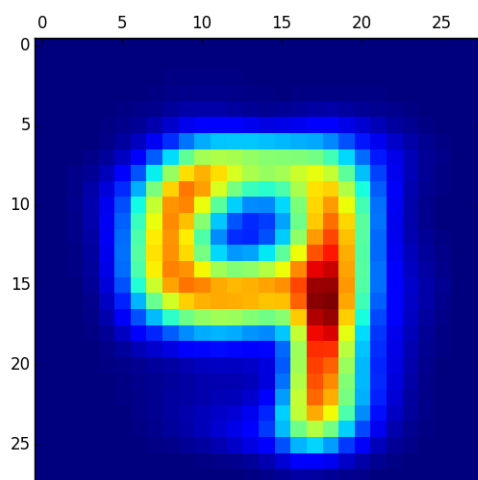
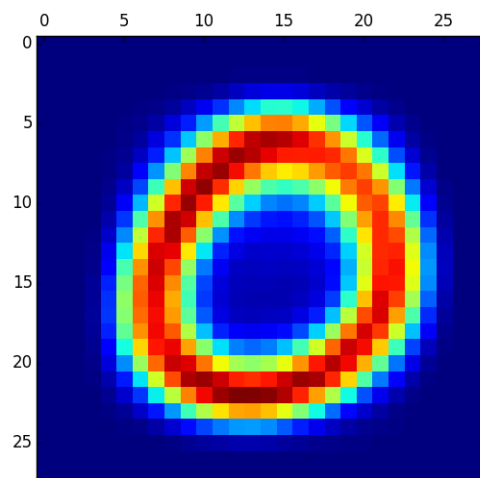
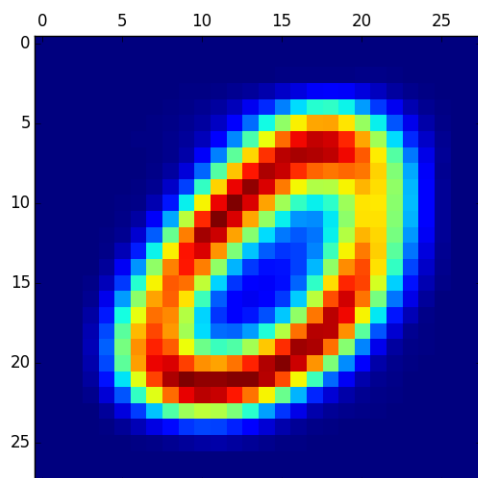
1.2 . With $K = 5$, visualize the cluster centers as below.

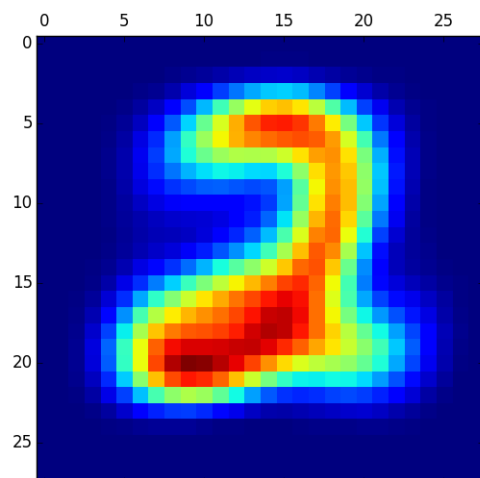
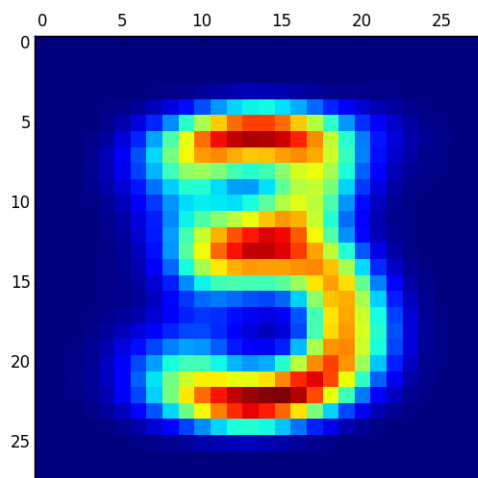




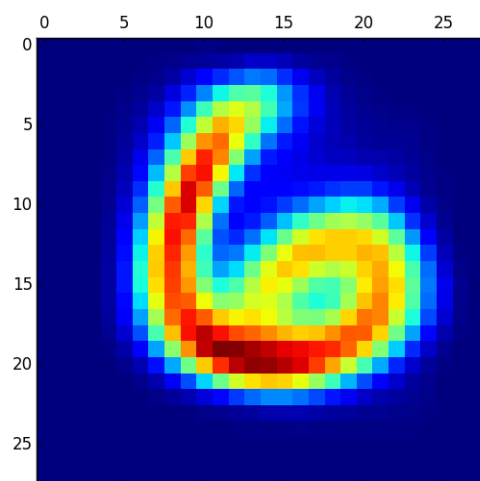
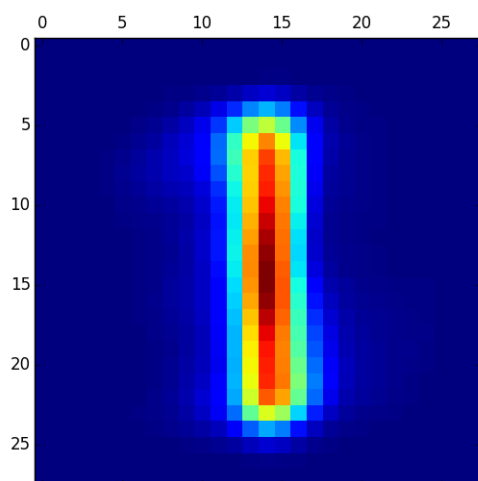
With $K = 10$, visualize the cluster centers as below.

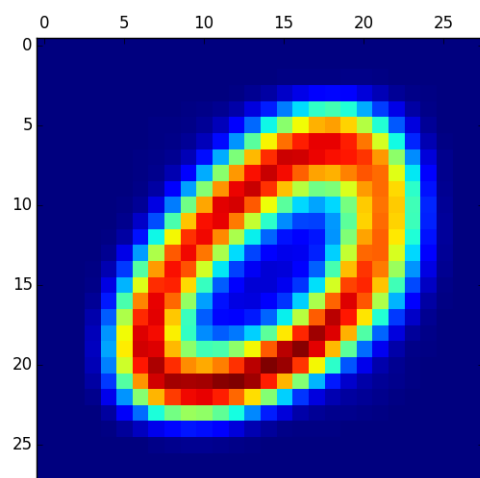
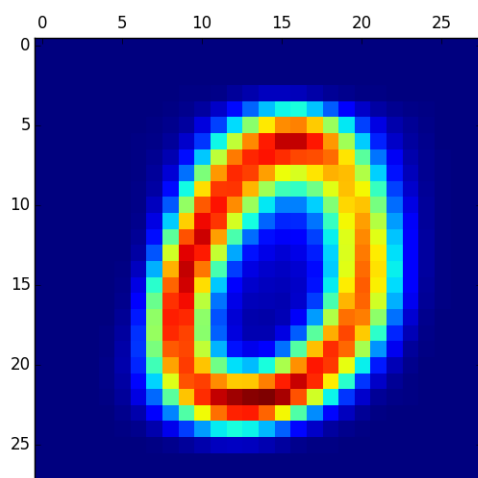
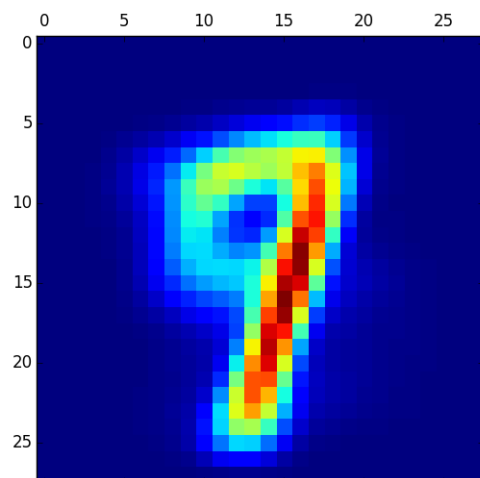
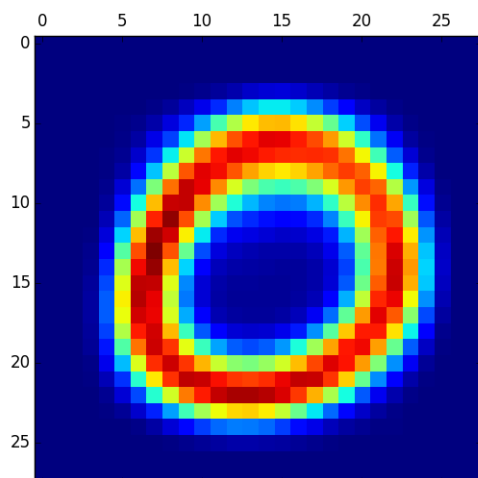


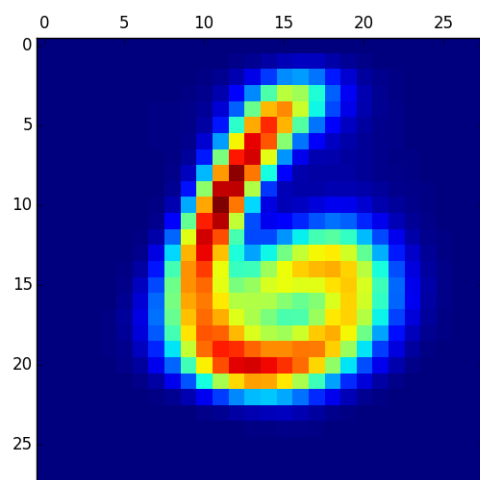
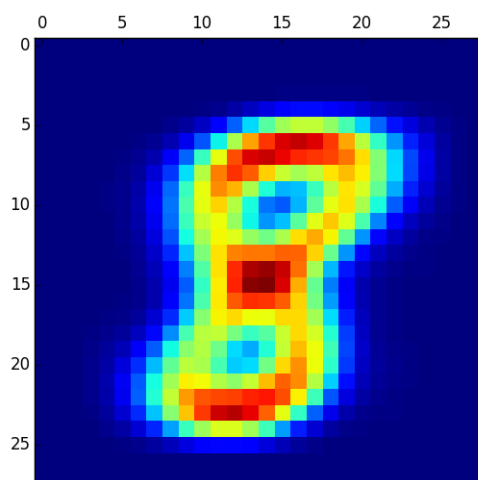
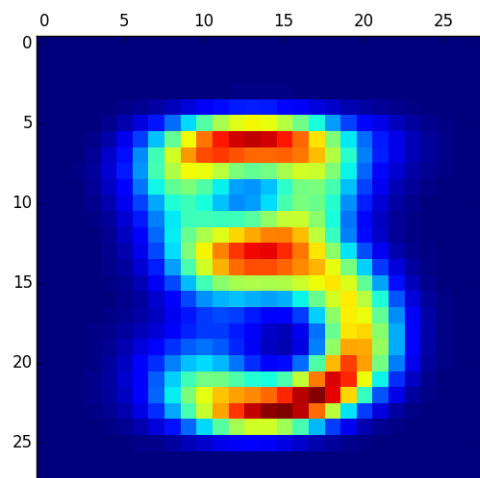
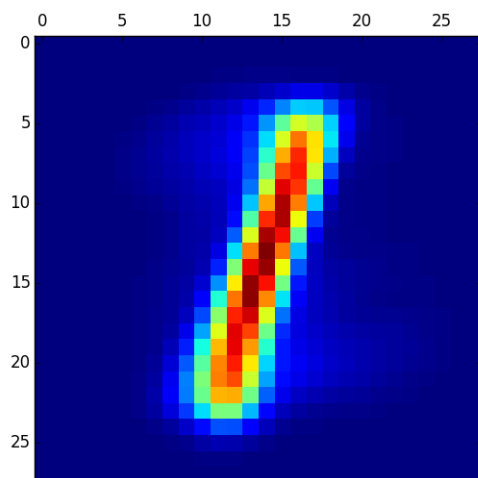


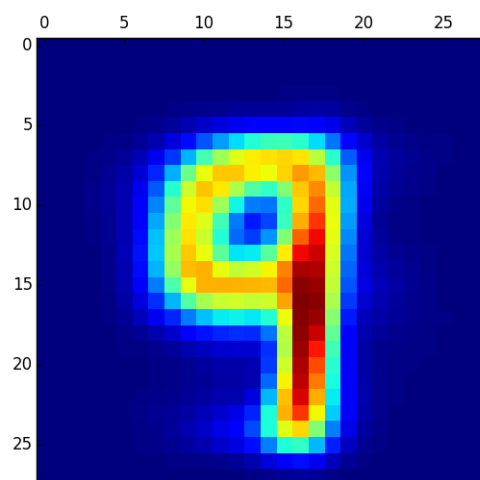
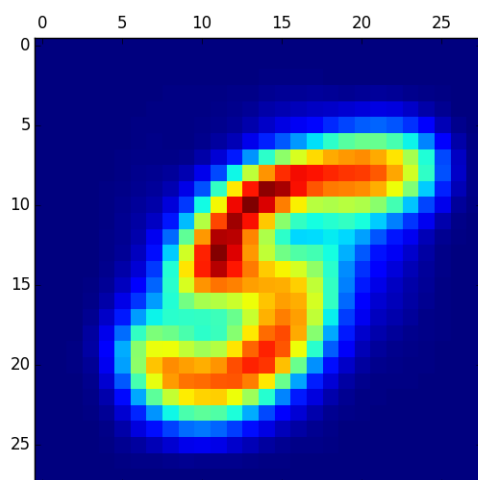
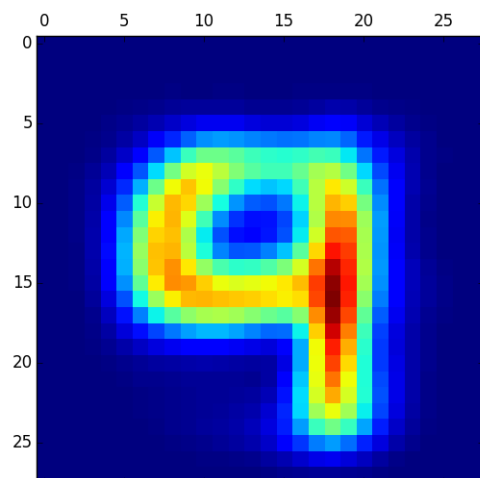
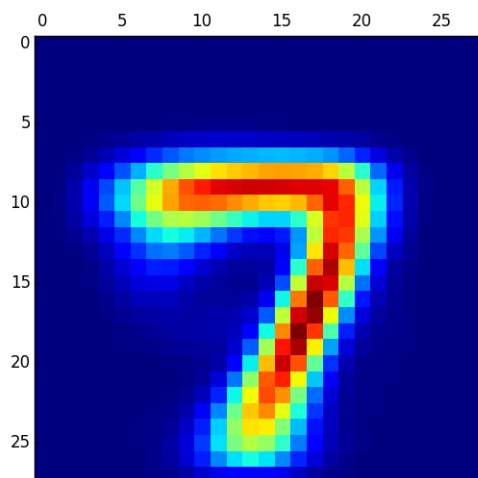


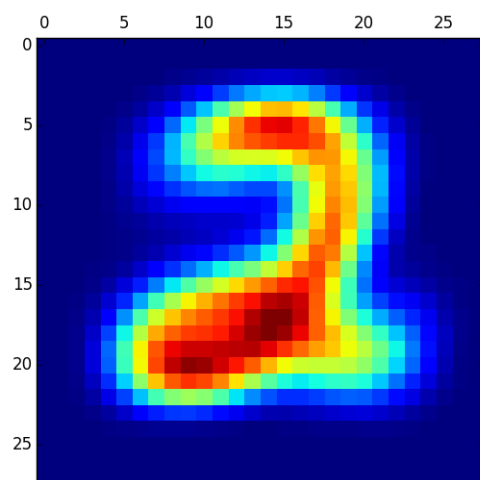
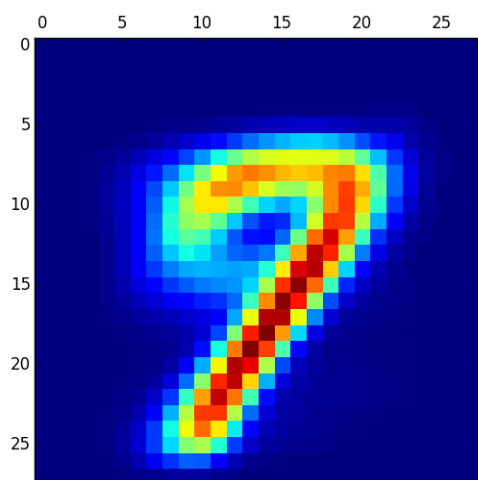
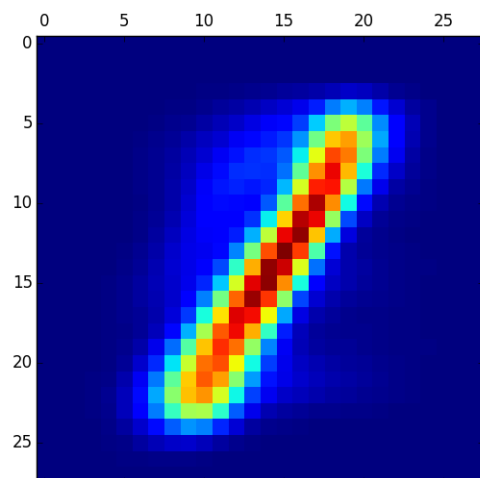
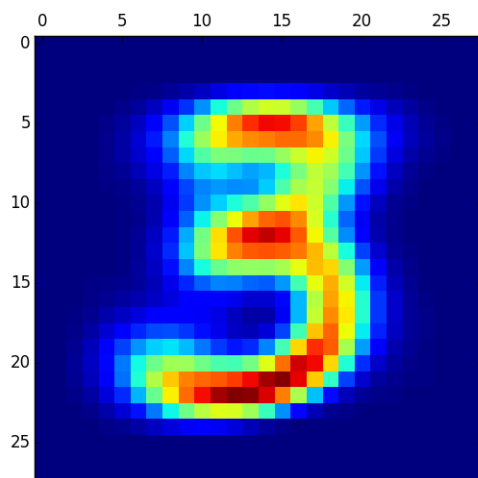
With $K = 20$, visualize the cluster centers as below.

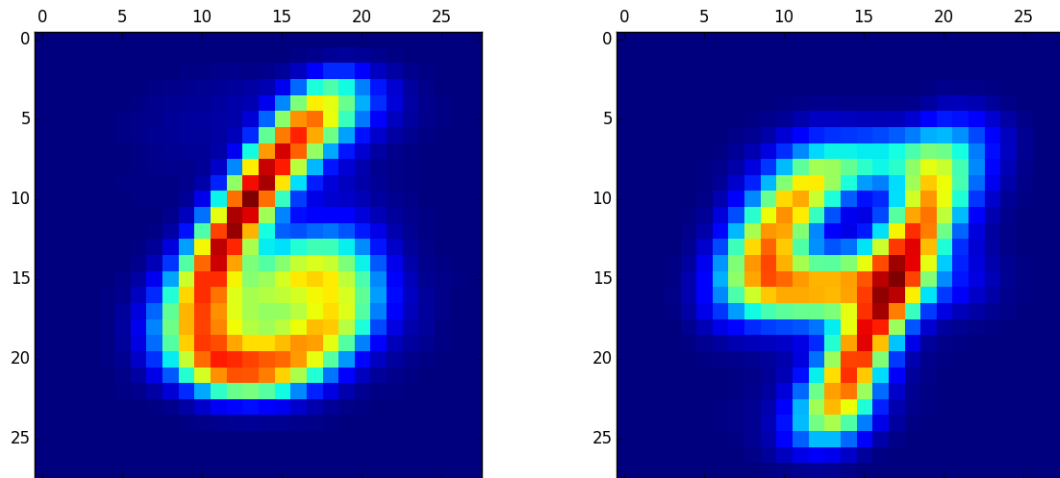












1.3 .

the k-mean loss does vary in different runs.

Problem 2

2.2 .

Recommend the joke by its average rating in the training set, the validation set accuracy is 0.6203252032520326

Find the k nearest neighbors of him/her, then make the prediction by averaging the ratings of these neighbors.

With $K = 10$, the validation set accuracy is 0.6490514905149052

With $K = 100$, the validation set accuracy is 0.6894308943089431

With $K = 1000$, the validation set accuracy is 0.6940379403794038

The accuracies are all higher than the accuracy that we got from the simple system.

2.3.1 . See code in q2.3.py

2.3.2 .

MSE is 20257100.7948 for $d = 2$

Validation Set accuracy is 0.7094850948509485 for $d = 2$

MSE is 18839053.5157 for $d = 5$

Validation Set accuracy is 0.713550135501355 for $d = 5$

MSE is 17092099.6887 for $d = 10$

Validation Set accuracy is 0.7132791327913279 for $d = 10$

MSE is 14192908.8435 for $d = 20$

Validation Set accuracy is 0.6859078590785908 for $d = 20$

Thus, MSE is decrease as variable d increase.

2.3.3 . See code in q2.3.3.py

2.3.4 .

MSE is 15111947.0226 for $d = 2$

Validation Set accuracy is 0.7073170731707317 for $d = 2$

MSE is 12380230.026 for $d = 5$

Validation Set accuracy is 0.7092140921409215 for $d = 5$

MSE is 9632977.97809 for $d = 10$

Validation Set accuracy is 0.7184281842818429 for $d = 10$

MSE is 5656292.52527 for $d = 20$

Validation Set accuracy is 0.6658536585365854 for $d = 20$

Compare to step 2, for the same d value, the MSE in is much lower in step 3. The validation set accuracies are almost same in both step 2 and 3, not having significant difference.

2.4 .

The best kaggle result is 0.72573

Code Q1

```

1  import scipy.io as sio
2  import numpy as np
3  from numpy.linalg import inv
4  import matplotlib.pyplot as plt
5  from sklearn.utils import shuffle
6
7
8  MnistData = sio.loadmat('./mnist_data/images.mat')
9  Images = MnistData['images']
10
11
12  x_dim = len(Images)
13  y_dim = len(Images[0])
14  image_index = len(Images[0][0])
15  Images = Images.transpose((2, 0, 1))
16  Images = Images.reshape(image_index, x_dim * y_dim)
17
18  def Kmean(images, k):
19      #initial clusters and centers
20      centers = np.random.randn(k, images[0].size)
21      clusters = [[] for i in range(k)]
22      for point in images:
23          sqaures = [np.square(np.linalg.norm(center - point)) for center in centers]
24          index = np.argmin(sqaures)
25          clusters[index] += [point]
26      centers = [np.mean(data, axis=0) for data in clusters]
27
28      #repeat assign point to new centers until no change
29      changed = True
30      iteration = 0
31      while changed:
32          print(iteration)
33          changed = False
34          newclusters = [[] for i in range(k)]
35          for i in range(k):
36              cluster = clusters[i]
37              for point in cluster:
38                  sqaures = [np.square(np.linalg.norm(center - point)) for center in centers]
39                  index = np.argmin(sqaures)
40                  newclusters[index] += [point]
41                  if i != index:
42                      changed = True
43          clusters = newclusters
44          centers = [np.mean(data, axis=0)for data in clusters]
45          iteration += 1
46      return centers
47
48  k = 5

```

```
49 centers = Kmean(Images, k)
50 print("done for k = 5")
51 for i in range(k):
52     plt.matshow(np.reshape(centers[i], (28, 28)))
53     plt.show()
54
55 k = 10
56 centers = Kmean(Images, k)
57 print("done for k = 10")
58 for i in range(k):
59     plt.matshow(np.reshape(centers[i], (28, 28)))
60     plt.show()
61
62 k = 20
63 centers = Kmean(Images, k)
64 print("done for k = 20")
65 for i in range(k):
66     plt.matshow(np.reshape(centers[i], (28, 28)))
67     plt.show()
```

Code Q2.2

```

1  import scipy.io as sio
2  import numpy as np
3  from numpy.linalg import inv
4  import matplotlib.pyplot as plt
5  from sklearn.utils import shuffle
6  from sklearn.preprocessing import Imputer
7  import operator
8
9  JokeData = sio.loadmat('./joke_data/joke_train.mat')
10 Images = JokeData['train']
11
12 data = [[int(i) for i in (line.strip().split(',') if line in open("./joke_data/validation.txt", 'r'))]]
13 data = np.array(data)
14 ValidationSet = data[:, :2]
15 ValidationLabels = data[:, 2]
16
17 def calculateAccuracy(expect, actual):
18     same = [i for i in range(len(actual)) if expect[i] == actual[i]]
19     return len(same) / len(actual)
20
21 #get nearest Neighbors with increase in distance
22 def findNearestNeighbors(Images, k, point):
23     distances = []
24     for i in range(len(Images)):
25         dist = np.linalg.norm(point - Images[i])
26         distances.append((Images[i], dist))
27     distances = sorted(distances, key=operator.itemgetter(1))
28     neighbors = []
29     for i in range(1, k + 1):
30         neighbors.append(distances[i][0])
31     return np.array(neighbors)
32
33 #2.2 Warm-up predict with average rating
34 averageRateValue = np.nanmean(Images, axis=0)
35 averageRatePredict = [ 1 if value > 0 else 0 for value in averageRateValue]
36 predict = []
37 for data in ValidationSet:
38     index = data[1]
39     predict.append(averageRatePredict[index - 1])
40
41 accuracy = calculateAccuracy(predict, ValidationLabels)
42 print("Rating by its average rating, the validation set accuracy is", accuracy)
43
44 #2.2 Warm-up with predict the k nearest neighbors
45
46 #replace Nan with 0
47 Ks = [10, 100, 1000]
48 Images[np.isnan(Images)] = 0

```

```
49 accuracy = 0
50 for k in Ks:
51     print('k =', k)
52     predict = []
53     userId = 0
54     averageRatePredict = []
55     for data in ValidationSet:
56         newuserId = data[0]
57         if not (newuserId == userId):
58             NearestNeighbors = findNearestNeighbors(Images, k, Images[newuserId - 1])
59             averageRateValue = np.mean(NearestNeighbors, axis=0)
60             averageRatePredict = [1 if value > 0 else 0 for value in averageRateValue]
61             predict.append(averageRatePredict[data[1] - 1])
62             userId = newuserId
63         else:
64             predict.append(averageRatePredict[data[1] - 1])
65     accuracy = calculateAccuracy(predict, ValidationLabels)
66     print('the validation set accuracy is ', accuracy)
67
68
```

Code Q2.3.2

```

1  import scipy.io as sio
2  import numpy as np
3  from sklearn import preprocessing
4  from sklearn.preprocessing import Imputer
5  import operator
6  import csv
7
8  JokeData = sio.loadmat('./joke_data/joke_train.mat')
9  Images = JokeData['train']
10
11 data = [[int(i) for i in (line.strip().split(',')')] for line in open("./joke_data/validation.txt", 'r')]
12 data = np.array(data)
13 ValidationSet = data[:, :2]
14 ValidationLabels = data[:, 2]
15
16 query = [[int(i) for i in (line.strip().split(',')')] for line in open("./joke_data/query.txt", 'r')]
17 query = np.array(query)
18 query = query[:, 1:3]
19
20 def calculateAccuracy(expect, actual):
21     same = [i for i in range(len(actual)) if expect[i] == actual[i]]
22     return len(same) / len(actual)
23
24 def predictWithUV(Dataset, vectorU, vectorV):
25     predict = []
26     for data in Dataset:
27         estimate = vectorU[data[0] - 1].dot(vectorV[:, data[1] - 1])
28         predict += [1] if estimate > 0 else [0]
29     return predict
30
31 # newImages replace Nan with 0
32 newImages = Images
33 newImages[np.isnan(newImages)] = 0
34 # newImages = preprocessing.normalize(newImages.astype("float"), norm='l2', axis=0)
35
36 # 2.3.1 compute U and V
37 U, s, V = np.linalg.svd(newImages, full_matrices=False)
38 U = U.dot(np.diag(s))
39
40 #2.3.2 with d vary compute MSE and validation set accuracies
41 Ds = [2, 5, 10, 20]
42 for d in Ds:
43     newU = U[:, :d+1]
44     newV = V[:, d+1, :]
45     newR = newU.dot(newV)
46     MSE = np.square(np.linalg.norm(newR - newImages))
47     print('MSE is', MSE, 'for d = ', d)
48     predict = predictWithUV(ValidationSet, newU, newV)

```



```

49     accuracy = calculateAccuracy(predict, ValidationLabels)
50     print('Validation Set accuracy is ', accuracy, 'for d = ', d)
51
52     # creat Kaggle submission file
53     predict_labels = predictWithUV(query, newU, newV)
54     indexes = [i for i in range(1, len(query) + 1)]
55     data = []
56     data += [indexes]
57     data += [predict_labels]
58     data = np.transpose(data, (1, 0)).tolist()
59     first_row = [['Id', 'Category']]
60     with open('2spampredict d = ' + str(d) + '.csv', 'w') as f:
61         a = csv.writer(f)
62         a.writerow(first_row)
63         a.writerows(data)
64

```

Code Q2.3.3

```

1  import scipy.io as sio
2  import numpy as np
3  from sklearn import preprocessing
4  from sklearn.preprocessing import Imputer
5  import operator
6  import csv
7
8  JokeData = sio.loadmat('./joke_data/joke_train.mat')
9  Images = JokeData['train']
10
11 data = [[int(i) for i in (line.strip().split(',') for line in open("./joke_data/validation.txt", 'r'))]]
12 data = np.array(data)
13 ValidationSet = data[:, :2]
14 ValidationLabels = data[:, 2]
15
16 query = [[int(i) for i in (line.strip().split(',') for line in open("./joke_data/query.txt", 'r'))]]
17 query = np.array(query)
18 query = query[:, 1:3]
19
20 #2.3.3
21 def calculateAccuracy(expect, actual):
22     same = [i for i in range(len(actual)) if expect[i] == actual[i]]
23     return len(same) / len(actual)
24
25 def predictWithUV(Dataset, vectorU, vectorV):
26     predict = []
27     for data in Dataset:
28         estimate = vectorU[data[0] - 1].dot(vectorV[:, data[1] - 1])
29         predict += [1] if estimate > 0 else [0]
30     return predict
31
32 def calculateLoss(Images, U, V, lamb):
33     loss = 0
34     for i in range(len(Images)):
35         for j in range(len(Images[0])):
36             if not np.isnan(Images[i][j]):
37                 loss += np.square(np.linalg.norm(U[i].T.dot(V[ : , j]) - Images[i][j]))\
38                     + lamb * np.square(np.linalg.norm(U[i]))\
39                     + lamb * np.square(np.linalg.norm(V[ : , j]))
40     return loss
41
42 def findUV(Images, d, lamb):
43     U = np.random.randn(len(Images), d)
44     V = np.random.randn(d, len(Images[0]))
45     # Changed = True
46     iterations = 0
47     while iterations < 100:
48         # Changed = False

```

```

49         iterations += 1
50         newU = []
51         origLoss = calculateLoss(Images, U, V, lamb)
52         for i in range(len(Images)):
53             left = lamb * np.identity(d)
54             right = np.zeros(d)
55             for j in range(len(Images[0])):
56                 if not np.isnan(Images[i][j]):
57                     left += np.outer(V[:, j], V[:, j])
58                     right += int(Images[i][j]) * V[:, j].T
59             newU.append(right.dot(np.linalg.inv(left)))
60         newU = np.array(newU)
61         U = newU
62
63         newV = []
64         for j in range(len(Images[0])):
65             left = lamb * np.identity(d)
66             right = np.zeros(d)
67             for i in range(len(Images)):
68                 if not np.isnan(Images[i][j]):
69                     left += np.outer(U[i].T, U[i].T)
70                     right += int(Images[i][j]) * U[i].T
71             newV.append(right.dot(np.linalg.inv(left)))
72         newV = np.array(newV).T
73         V = newV
74
75         newLoss = calculateLoss(Images, U, V, lamb)
76         # diff = origLoss - newLoss
77         # print('diff is', diff)
78         print('MSE is ', newLoss)
79         predict = predictWithUV(ValidationSet, U, V)
80         accuracy = calculateAccuracy(predict, ValidationLabels)
81         print('Validation Set accuracy is ', accuracy, 'for d = ', d)
82         # if abs(diff) > 3:
83         #     Changed = True
84     # creat Kaggle submission file
85     predict_labels = predictWithUV(query, U, V)
86     indexes = [i for i in range(1, len(query) + 1)]
87     data = []
88     data += [indexes]
89     data += [predict_labels]
90     data = np.transpose(data, (1, 0)).tolist()
91     first_row = [['Id', 'Category']]
92     with open('3spampredict d = ' + str(d) + '.csv', 'w') as f:
93         a = csv.writer(f)
94         a.writerow(first_row)
95         a.writerows(data)
96     return U, V
97
98

```

```
99  lamb = 0.1
100  Ds = [2, 5, 10, 20]
101
102  for d in Ds:
103      print('begin train with d = ', d)
104      U, V = findUV(Images, d, lamb)
105      MSE = calculateLoss(Images, U, V, lamb)
106      print('MSE is', MSE, 'for d = ', d)
107      predict = predictWithUV(ValidationSet, U, V)
108      accuracy = calculateAccuracy(predict, ValidationLabels)
109      print('Validation Set accuracy is ', accuracy, 'for d = ', d)
110      print('done train with d = ', d)
111
112
```