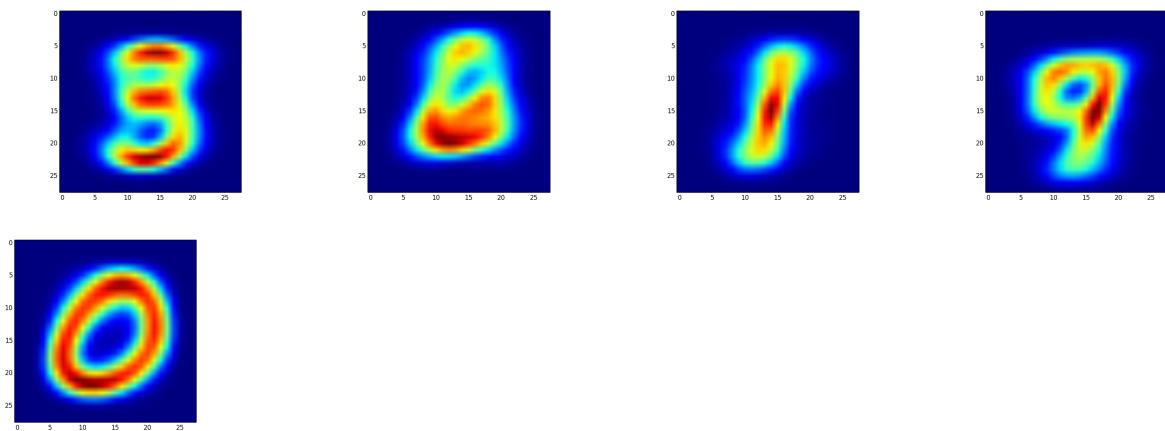


# CS189—Spring 2016 — Homework 7Solutions

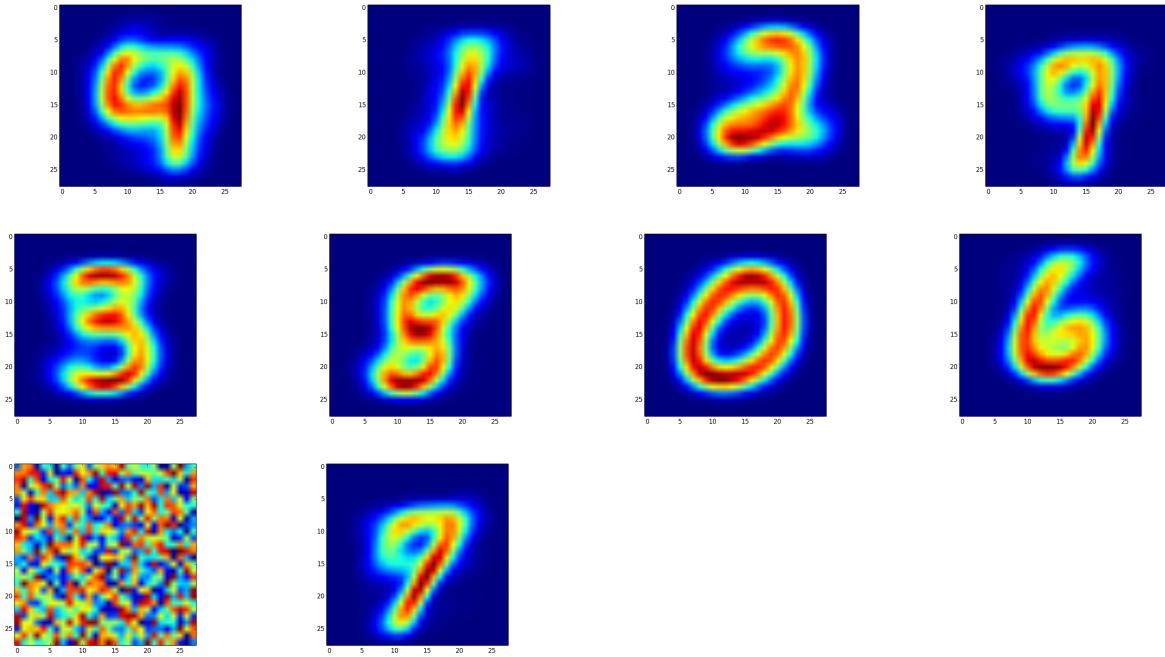
Andy Chu, SID 23042499, PUT SOMETHING HERE

# 1 Problem

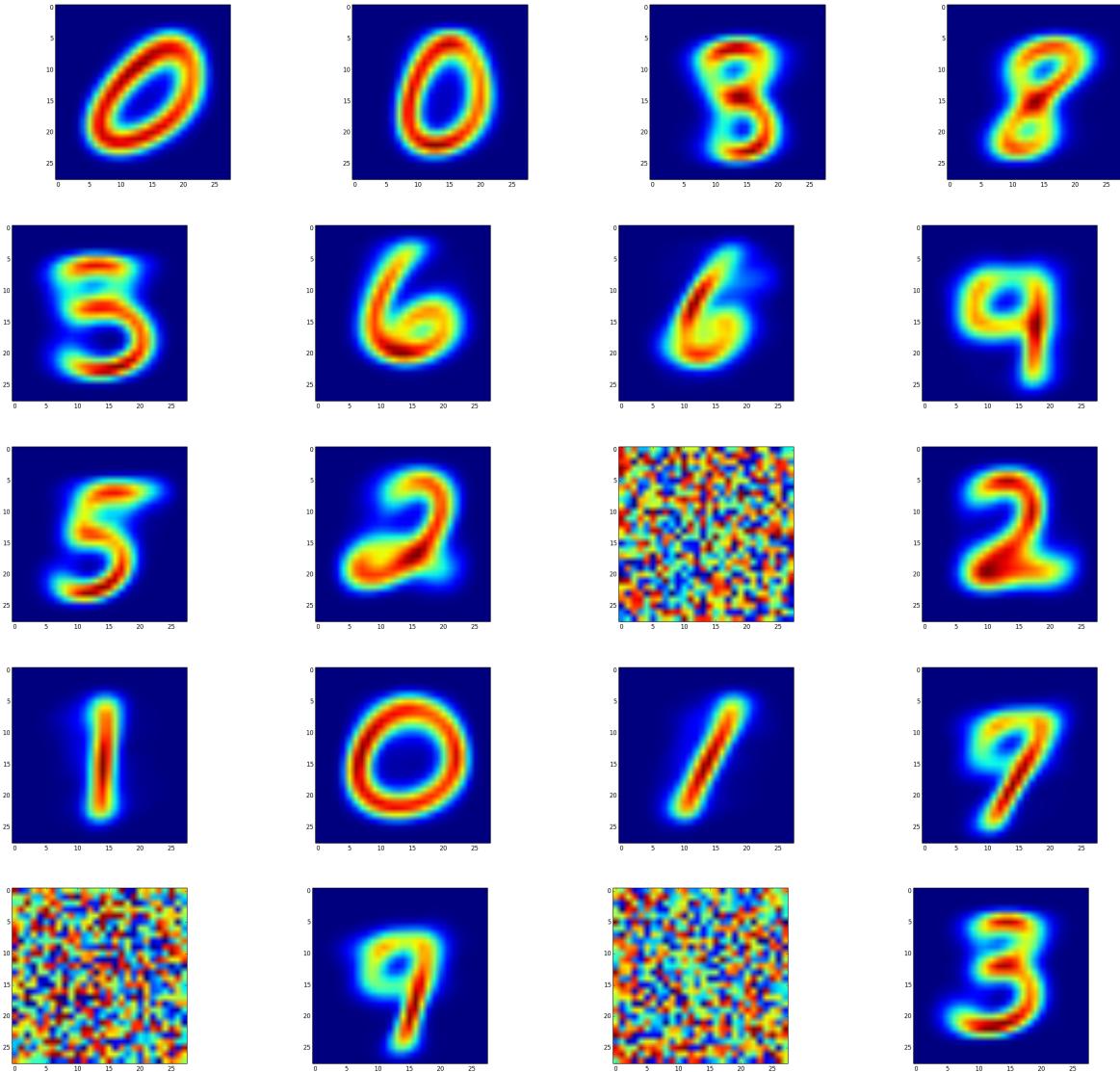
$k = 5$  for 50 iterations:



$k = 10$  for 50 iterations:



$k = 20$  for 50 iterations:



K-mean loss does vary in different runs. When  $k = 5$  was run for 10 iterations, the loss was at first 102,135,443.731 and then it was 103,592,741.946 when run again.

## 2 Problem 2.2

For the simplest recommender system, the accuracy on validation set is 0.6203. For the advanced method with  $k = 10$ , the accuracy is 0.6491. The advanced method with  $k = 100$ , the accuracy is 0.6894. The advanced method with  $k = 1,000$ , the accuracy is 0.6940. The accuracies of the advanced system are all better than the accuracy of the simple system.

### 3 Problem 2.3

Latent vector with mean square error for  $d = 2$  is 0.6477 with MSE of 20,703,188. Latent vector with mean square error for  $d = 5$  is 0.6667 with MSE of 19,192,013. Latent vector with mean square error for  $d = 10$  is 0.6520 with MSE of 17,398,454. Latent vector with mean square error for  $d = 20$  is 0.6284 with MSE of 14,458,503.

We see that as we increase  $d$ , the MSE decreases.

When we change our loss function to take into account sparse data, having  $k$  is 10 clusters, threshold of 0.1, learning rate for gradient descent of 0.0001, and lambda regularization 100, the MSE is 1,2818,257 and the accuracy is 0.7255. We see that the MSE with this loss function is smaller than that of MSE when loss function is mean square. The validation accuracy for loss function that takes into account sparse data is also better than the validation accuracy of loss function for mean square.

Kaggle score: 0.71434

## 4 Appendix

python kmeans.py to run kmeans for 1

```
1 import numpy as np
2 from scipy import io, spatial, misc
3 import scipy.misc
4 from matplotlib import pyplot as plt
5 import random
6
7
8
9
10
11
12
13
14 train_contents = io.loadmat('images.mat')
15
16 images = train_contents['images']
17
18
19 pixel_images = images[0,:,:]
20 for i in range(1, 28):
21
22     pixel_images = np.vstack([pixel_images, images[i,:,:]])
23
24
25
26 pixel_images = pixel_images.T
27
28
29 class KMeans(object):
30     def __init__(self, k=10, iters=300):
31         self.k = k
32         self.centroids = {}
33         self.iters = iters
34
35
36
37
38     def train(self, data):
39         #initialize the centroids randomly with 784 pixels with values up to 255
40         for i in range(self.k):
41             self.centroids[i] = np.random.choice(256, 784) #256 because the max
42             pixel in 255
43
44             iters = 0
45
46             while iters < self.iters:
47                 clusters = {}
48                 J = 0.0
49                 index = 0
50                 #if ((iters % 10) == 0):
51                 print "training", iters
52                 for s in data:
53                     minDist = 3000000      #initialize to number bigger than 255*784 =
54                     200,175
55                     label = -1
```

```

54         # find the label c that has smalled euclidean distance with data
55         point
56             for c in self.centroids:
57                 dist = spatial.distance.euclidean(s, self.centroids[c])
58                 if dist < minDist:
59                     minDist = dist
60                     label = c
61
62             #create new group if smallest label is not in clusters
63             if label not in clusters:
64                 clusters[label] = []
65
66             #add the data to cluster it belongs to
67             clusters[label].append(index)
68
69             index += 1
70             # take mean of each cluster to find the new centroid
71             # increment the loss to figure out new loss
72             for c in clusters:
73                 self.centroids[c] = np.mean(data[clusters[c]], axis=0)
74                 for g in data[clusters[c]]:
75                     J += spatial.distance.euclidean(g, self.centroids[c])
76
77             print "Loss is", J, "for iteration ", iters
78
79             iters += 1
80
81             print "final loss of", J, "for iteration ", iters
82 iters = 50
83
84 km5 = KMeans(5, iters)
85 km5.train(pixel_images)
86
87 for n in km5.centroids:
88     np.savetxt('q1/km5/' + str(n), km5.centroids[n].reshape(28,28))
89
90 #after write files, then can use this for loop to plot each of the images
91 for i in range(5):
92     p = np.loadtxt('q1/km5/' + str(i))
93
94     plt.imshow(p)
95     plt.savefig('q1/km5/' + str(i))
96
97
98 km10 = KMeans(10, iters)
99 km10.train(pixel_images)
100
101 for n in km10.centroids:
102     np.savetxt('q1/km10/' + str(n), km10.centroids[n].reshape(28,28))
103
104 #after write files, then can use this for loop to plot each of the images
105 for i in range(10):
106     p = np.loadtxt('q1/km10/' + str(i))
107
108     plt.imshow(p)
109     plt.savefig('q1/km10/' + str(i))
110
111

```

```
112 km20 = KMeans(20, iters)
113 km20.train(pixel_images)
114
115 for n in km20.centroids:
116     np.savetxt('q1/km20/' + str(n), km20.centroids[n].reshape(28,28))
117
118 #after write files, then can use this for loop to plot each of the images
119 for i in range(20):
120     p = np.loadtxt('q1/km20/' + str(i))
121
122     plt.imshow(p)
123     plt.savefig('q1/km20/' + str(i))
```

python avg.py to run simple recommender system and advanced system in 2.2

```

1 import numpy as np
2 from scipy import io, spatial, misc
3 import scipy.misc
4 from matplotlib import pyplot as plt
5 from scipy import io, spatial, misc
6 import random
7
8
9 train_contents = io.loadmat('joke_train.mat')
10
11 train = train_contents['train']
12
13 validation = 'validation.txt'
14
15 class recommend_avg(object):
16     def __init__(self, data):
17         self.avg = np.array([np.nanmean(data, axis=0)]).T
18         #print self.avg.shape
19
20     def predict(self, file):
21         pred = []
22         f = open(file, 'r')
23         for line in f:
24             line = line.split(',')
25             joke_num = int(line[1]) - 1 # - 1 because of how python starts with 0 for
26             index
27             joke_avg_rating = self.avg[joke_num][0]
28
29             if joke_avg_rating > 0:
30                 pred.append(1)
31             else:
32                 pred.append(0)
33         return pred
34
35 def accuracy(prediction, validation):
36     correct = 0.0
37     total = 0
38     f = open(validation, 'r')
39     for line in f:
40         line = line.split(',')
41         if (int(line[2])) == int(prediction[total]):
42             correct += 1
43         total += 1
44     return correct/float(total)
45
46 ra = recommend_avg(train)
47 prediction = ra.predict(validation)
48 print "standard accuracy is", accuracy(prediction, validation)
49
50
51 class advanced(object):
52     def __init__(self, data, k):
53         #self.data = np.array([np.nan_to_num(data)])
54         self.data = np.nan_to_num(data) #(24983, 100)
55         self.k = k
56         self.avg = {}

```



```
111     print "on", i
112     self.avg = avgcluster
113     #print len(self.avg), "shape of avg cluster"
114
115
116     def predict(self, file):
117         pred = []
118         f = open(file, 'r')
119         for line in f:
120             line = line.split(',')
121             joke_num = int(line[1]) - 1 # - 1 because of how python starts with 0 for
122             index
123             user = int(line[0]) - 1
124             #print "self.avg[user]", self.avg[user]
125             #print "joke num", joke_num
126             #print "value is", self.avg[user][joke_num]
127             #print "user", user, "joke num", joke_num
128             joke_avg_rating = self.avg[user][joke_num]
129
130             if joke_avg_rating > 0:
131                 pred.append(1)
132             else:
133                 pred.append(0)
134         return pred
135
136
137     adv = advanced(train, 10)
138     #print "adv object", adv.data, adv.k, adv.avg
139     adv.train()
140     prediction = adv.predict(validation)
141     print "accuracy advanced 10", accuracy(prediction, validation)
142
143     adv = advanced(train, 100)
144     #print "adv object", adv.data, adv.k, adv.avg
145     adv.train()
146     prediction = adv.predict(validation)
147     print "accuracy advanced 100", accuracy(prediction, validation)
148
149     adv = advanced(train, 1000)
150     #print "adv object", adv.data, adv.k, adv.avg
151     adv.train()
152     prediction = adv.predict(validation)
153     print "accuracy advanced 1000", accuracy(prediction, validation)
```

python latent.py to run latent factor model with mean square loss function and loss function that takes into account sparse data for 2.3

```

1 import numpy as np
2 from scipy import io, spatial, misc
3 import scipy.misc
4 from matplotlib import pyplot as plt
5 import random
6
7 DEBUG = 0
8
9
10 #print random.__file__
11
12 #after write files, then can use this for loop to plot each of the images
13
14
15
16 train_contents = io.loadmat('joke_train.mat')
17
18 images = train_contents['train'] #(24983,100)
19
20 orig = images
21
22 validation = 'validation.txt'
23
24 testing = 'query.txt'
25
26 class joke(object):
27     def __init__(self, d):
28         self.d = d
29         self.u = []          #(d, 100)
30         self.v = []          #(24983, d)
31
32     def train(self, data):
33         image = np.nan_to_num(data)
34         mean = np.mean(image, axis=0)
35         image = image - mean
36         U, s, V = np.linalg.svd(image, full_matrices=False)
37         sort = np.argsort(s)[::-1]           #[::-1] <-- to reverse the order
38         #print "image", image, "mean", mean, "new image", image
39         #print U.shape, V.shape, s.shape
40         #print "sort", sort, "np.argsort", np.argsort(s)
41         d = self.d
42         U = U[:, sort]
43         V = V.T[:, sort]
44         S = np.diag(s[sort])
45         self.u = np.dot(U[:, :d], np.sqrt(S[:d, :d]))
46         self.v = np.dot(np.sqrt(S[:d, :d]), V[:, :d].T)
47         #print "v", self.v.shape, "u", self.u.shape
48         print "MSE for " + str(self.d)
49         mse = np.nansum((np.dot(self.u, self.v) - image)**2)
50         print mse
51
52     def predict(self, file):
53         pred = []
54         f = open(file, 'r')
55         for line in f:
56             line = line.split(',')

```

```

57     user = int(line[0]) - 1
58     joke_num = int(line[1]) - 1 # - 1 because of how python starts with 0 for
      index
59     joke_avg_rating = np.dot(self.u[user,:], self.v[:,joke_num])
60
61     if joke_avg_rating > 0:
62         pred.append(1)
63     else:
64         pred.append(0)
65     return pred
66
67 def accuracy(prediction, validation):
68     correct = 0.0
69     total = 0
70     f = open(validation, 'r')
71     for line in f:
72         line = line.split(',')
73         if (int(line[2])) == int(prediction[total]):
74             correct += 1
75         total += 1
76     return correct/float(total)
77
78
79 j = joke(2)
80 j.train(images)
81 prediction = j.predict(validation)
82 print "accuracy for meansq " + str(2), accuracy(prediction, validation)
83
84 j = joke(5)
85 j.train(images)
86 prediction = j.predict(validation)
87 print "accuracy for meansq " + str(5), accuracy(prediction, validation)
88
89 j = joke(10)
90 j.train(images)
91 prediction = j.predict(validation)
92 print "accuracy for meansq " + str(10), accuracy(prediction, validation)
93
94 j = joke(20)
95 j.train(images)
96 prediction = j.predict(validation)
97 print "accuracy for meansq" + str(20), accuracy(prediction, validation)
98
99 class adv(object):
100     def __init__(self, d, threshold, grad, lamb):
101         self.d = d
102         self.u = []          #(d, 100)
103         self.v = []          #(24983, d)
104         self.threshold = threshold
105         self.grad = grad
106         self.lamb = lamb
107     def train(self, data):
108         image = np.nan_to_num(data)
109         mean = np.mean(image, axis=0)
110         image = image - mean
111         U, s, V = np.linalg.svd(image, full_matrices=False)
112         sort = np.argsort(s)[::-1]           #[::-1] <-- to reverse the order
113         # print "image", image, "mean", mean, "new image", image
114         # print U.shape, V.shape, s.shape

```

```

115 # print "sort", sort, "np.argsort", np.argsort(s)
116 d = self.d
117 U = U[:, sort]
118 V = V.T[:, sort]
119 S = np.diag(s[sort])
120 self.u = np.dot(U[:, :d], np.sqrt(S[:d, :d]))
121 self.v = np.dot(np.sqrt(S[:d, :d]), V[:, :d].T)
122 # print "v", self.v.shape, "u", self.u.shape
123 # print "loss for " + str(self.d)
124 #print "dot of u and v", np.dot(self.u, self.v).shape
125 print (np.dot(self.u, self.v) - orig)**2
126 print np.sum(np.multiply(self.u, self.u))
127 print np.sum(np.multiply(self.v, self.v))
128 loss = np.nansum((np.dot(self.u, self.v) - orig)**2) + self.lamb* (np.sum(np.
129 multiply(self.u, self.u)) + np.sum(np.multiply(self.v, self.v)))
130 print "first loss", loss
131 iteration = 0
132 prev = loss
133 while True:
134     #print "u", self.u, "u shape", self.u.shape
135     #print "v", self.v, "v shape", self.v.shape
136     #print "dot product", np.dot(self.u, self.v), "dot shape", np.dot(self.u, self.
137     .v).shape
138     #print "image", image, "image shape", image.shape
139     #print "subtract image from dot product", np.dot(self.u, self.v) - image
140     #self.u -= 2*(np.dot(np.nan_to_num((np.dot(self.u, self.v) - orig)), self.v.T)
141     + self.lambd*self.u)
142     if (iteration == 501):
143         print "converged on iteration ", iteration, " with loss of ", loss
144         return
145
146     self.u -= 2 * self.grad*(np.dot((np.nan_to_num(np.dot(self.u, self.v) - orig)),
147                                     self.v.T) + self.lamb * self.u)
148     #self.v -= 2*(np.dot(self.u.T, np.nan_to_num(np.dot(self.u, self.v) - orig)) +
149     self.lambd*self.v)
150     self.v -= 2 * self.grad*(np.dot(self.u.T, np.nan_to_num((np.dot(self.u, self.v) -
151                                     orig))) + self.lamb * self.v)
152     #print "new u", self.u
153     #print "new v", self.v
154     loss = np.nansum((np.dot(self.u, self.v) - orig)**2) + self.lamb* (np.sum(np.
155     multiply(self.u, self.u)) + np.sum(np.multiply(self.v, self.v)))
156     if (np.abs(prev - loss) < self.threshold):
157         print "converged on iteration ", iteration, " with loss of ", loss
158         return
159     if ((iteration % 100) == 0):
160         print "iteration ", iteration
161         print "loss of ", loss
162         print "prev - loss", np.abs(prev - loss)
163         prediction = j.predict(validation)
164         print "accuracy for", iteration, "is", accuracy(prediction, validation)
165
166         prev = loss
167         iteration += 1
168
169 def predict(self, file):
170     pred = []
171     f = open(file, 'r')
172     for line in f:
173

```

```
167     line = line.split(',')
168     user = int(line[0]) - 1
169     joke_num = int(line[1]) - 1 # - 1 because of how python starts with 0 for
170     index
171     joke_avg_rating = np.dot(self.u[user,:], self.v[:,joke_num])
172
173     if joke_avg_rating > 0:
174         pred.append(1)
175     else:
176         pred.append(0)
177     return pred
178
179 def test(self, file):
180     pred = []
181     f = open(file, 'r')
182     w = open('kaggle_submission.txt', 'w')
183     w.write('Id,Category\n')
184     for line in f:
185         line = line.split(',')
186         user = int(line[1]) - 1
187         joke_num = int(line[2]) - 1 # - 1 because of how python starts with 0 for
188         index
189         joke_avg_rating = np.dot(self.u[user,:], self.v[:,joke_num])
190
191         if joke_avg_rating > 0:
192             w.write(str(line[0]) + ',', + str(1) + '\n')
193         else:
194             w.write(str(line[0]) + ',', + str(0) + '\n')
195     w.close()
196     return pred
197
198 j = adv(10, 0.1, 0.0001, 100)
199 j.train(images)
200 prediction = j.predict(validation)
201 print "accuracy for sparseAdjusted ", accuracy(prediction, validation)
202 j.test(testing)
```