

# EE378B Homework 8 Solution

Due to: Zixuan Liu

## Question (a)

We generated the synthetic data following the problem setting. The table shows that with  $n$  increasing, the ratio of average  $Q$  becomes larger. The code for part (a) and (b) is listed below.

Methods	n=200	n=400	n=800	n=1600	n=3200
PCA	0.905	0.946	0.971	0.986	0.993
NMF	0.909	0.949	0.975	0.987	0.992

Table 1: Results over 10 iterations

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from nmf_hw import nmf
4 import scipy as sp
5 from scipy.sparse.linalg import svds
6
7 d = 200
8 zeros = np.zeros(d // 2)
9 ones = np.ones(d // 2)
10 theta = np.concatenate([zeros, ones])
11
12 N = [200, 400, 800, 1600, 3200]
13 const = np.sqrt(d)
14 r = 1
15 Iter = 10
16 for n in N:
17     Q_avg_PCA = 0
18     Q_avg_NMF = 0
19     for iter in range(Iter):
20         np.random.shuffle(theta)
21         X = np.zeros((n, d))
22         w = np.random.rand(n)
23         for i in range(n):
24             for j in range(d):
25                 p = w[i] * theta[j] / const
26                 Xij = const * np.random.choice(2, p=[1-p, p])
27                 X[i, j] = Xij
28
29         # PCA
30         SIGMA = X.T.dot(X) / n
31         U, sigma, VT = sp.sparse.linalg.svds(SIGMA, k=r)
32         # U, sigma, VT = np.linalg.svd(SIGMA)
33         # print(sigma)
34         # print(U.shape)
35         u1 = U[:, 0]
36         u1_norm = np.linalg.norm(u1)
37         theta_norm = np.linalg.norm(theta)
38         Q = np.abs(np.dot(u1, theta)) / (u1_norm * theta_norm)
39         # print('n =', n, 'Iter', iter, ',Q = ', Q)
40         Q_avg_PCA += Q
41
42         # NMF
43         X_norm = X / (X.sum(axis=1, keepdims=True) + 1e-10)
44         Winit = np.random.rand(n, r) + 0.5
45         Hinit = np.random.rand(r, d) + 0.5
46         W, H = nmf(X_norm, Winit, Hinit, tol=1e-3, maxiter=100)
```

```

47 # print(H)
48 # print(theta)
49 h = H[0]
50 h_norm = np.linalg.norm(h)
51 theta_norm = np.linalg.norm(theta)
52 Q = np.abs(np.dot(h, theta)) / (h_norm * theta_norm)
53 # print('n =', n, 'Iter', iter, 'Q = ', Q)
54 Q_avg_NMF += Q
55 print('NMF n =', n, 'Q_avg = ', Q_avg_NMF / Iter)
56 print('PCA n =', n, 'Q_avg = ', Q_avg_PCA / Iter)

```

## Question (b)

We implemented NMF algorithm following the instruction of homework document. Since the name of algorithm is non-negative matrix factorization, we consider using non-negative matrix H and W as initialization. Also, we observe that the algorithm will gradually converge to some points, therefore we regard  $|F(X, H^{t+1}, W^{t+1}) - F(X, H^t, W^t)| < \text{tol}$  as the convergence criterion. Note that we choose  $10^{-3}$  for both synthetic data as well as MNIST test data.

The result of NMF shows that in most settings, it outperforms the PCA algorithm, except for the n=3200 case.

Methods	n=200	n=400	n=800	n=1600	n=3200
PCA	0.905	0.946	0.971	0.986	0.993
NMF	0.909	0.949	0.975	0.987	0.992

Table 2: Results over 10 iterations

```

1 ''' Nonnegative Matrix Facotrization'''
2 import numpy as np
3 from time import time
4
5 def nmf(X, Winit, Hinit, tol, maxiter=100):
6     eps_threshold = 1e-10
7     W = Winit # (n, r)
8     H = Hinit # (r, d)
9     initt = time()
10    past_score = 0
11    for iter in range(maxiter):
12        WH = W.dot(H) # (n, d)
13        X_WH = X / (WH + eps_threshold) # (n, d)
14        X_WH_H = X_WH.dot(H.T) # (n, r)
15        W_tilde = W * X_WH_H # (n, r)
16        newW = W_tilde / W_tilde.sum(axis=0, keepdims=True) # (n, r)
17        newWH = newW.dot(H) # (n, d)
18        X_newWH = X / (newWH + eps_threshold) # (n, d)
19        W_X_newWH = newW.T.dot(X_newWH) # (r, d)
20        H_tilde = H * W_X_newWH # (r, d)
21        newH = H_tilde / H_tilde.sum(axis=1, keepdims=True) # (r, d)
22        H = newH
23        W = newW
24        score = F(X, W, H)
25        if np.abs(score - past_score) < tol: break
26        past_score = score
27    return W, H
28
29 def F(X, W, H):
30     eps_threshold = 1e-10
31     n, d = X.shape
32     WH = W.dot(H) # (n, d)
33     score = np.sum(X * np.log(WH + eps_threshold) - WH)
34     return score

```

## Question (c)

We first extract the data, and apply PCA on it. From the plots of the 6 principal vectors, we can see numbers mixed up, for example, we can see a dark blue 6 from the first plot, and a dark blue 3 or 8 in the third plots, and a dark blue 9, a yellow 2 from the fourth plots. These implies that the PCA do extract some useful features from the dataset.

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import scipy as sp
4 from scipy.sparse.linalg import svds
5 from nmf_hw import nmf
6
7 fname = './mnist_test.csv'
8 f_outpath = './'
9 with open(fname, 'rb') as f:
10     data = np.loadtxt(f, delimiter=',')
11
12 print(data.shape)
13 print(data[:, 0])
14
15 y = data[:, 0]
16 x = data[:, 1:]
17 x_img = x.reshape(-1, 28, 28)
18
19 N, d = x.shape
20 r = 6
21 '''PCA for MNIST'''
22 SIGMA = x.T.dot(x) / N
23 U, sigma, VT = sp.sparse.linalg.svds(SIGMA, k=r)
24 for i in range(r):
25     u = U[:, i]
26     u_img = u.reshape(28, 28)
27     plt.imshow(u_img)
28     plt.title('PCA r=' + str(i + 1))
29     plt.colorbar(shrink=0.83)
30     # plt.show()
31     plt.savefig(f_outpath + 'PCA r= ' + str(i + 1) + '.png', format='png', transparent=True, dpi=300, pad_
32               inches = 0)
33     plt.clf()
34     plt.cla()
```

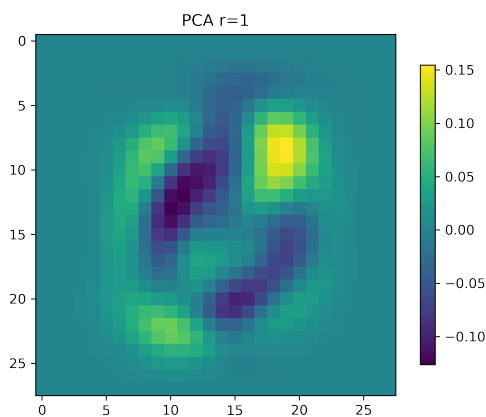


Figure 1: PCA for MNIST:  $r = 1$

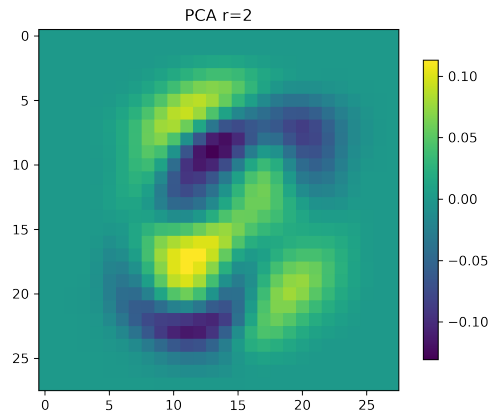


Figure 2: PCA for MNIST:  $r = 2$

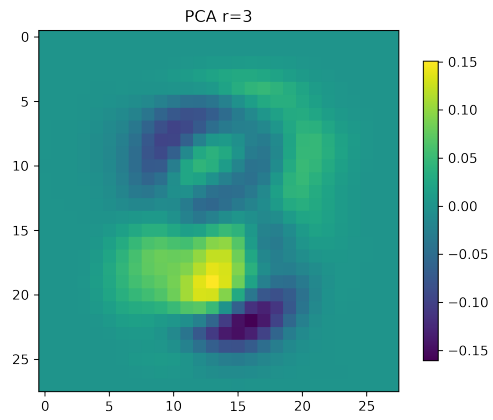


Figure 3: PCA for MNIST:  $r = 3$

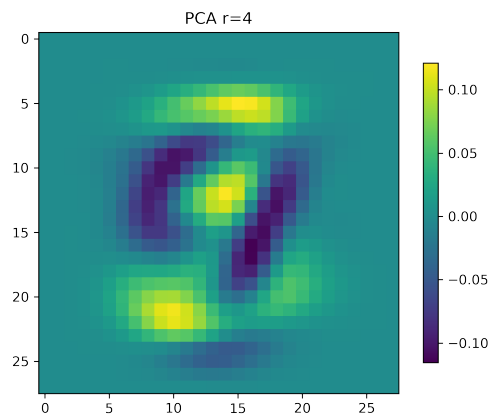


Figure 4: PCA for MNIST:  $r = 4$

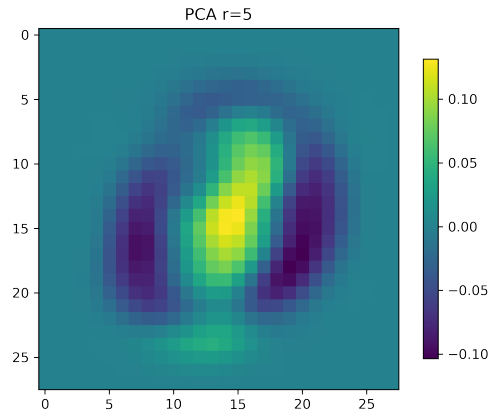


Figure 5: PCA for MNIST:  $r = 5$

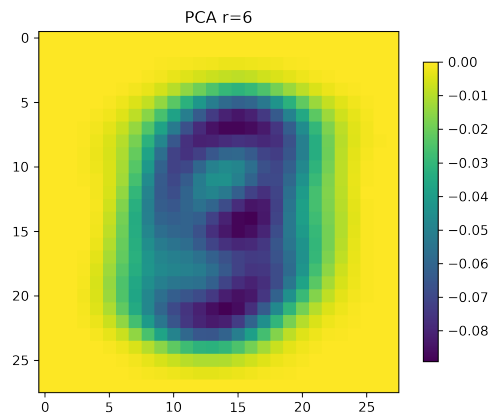


Figure 6: PCA for MNIST:  $r = 6$

## Question (d)

We apply the NMF algorithm to MNIST test data. Similar to PCA algorithm, we can see features from the plots. For example, we can see a 'C' like feature in the first plot, and a rough structure of 4 in the second, a line in the third plot, indicating the digit 1, 7, 9. We also see a rough 7 or 9 or 2 in the fourth plot, and a 9 in the fifth plot, and a 1 in the sixth plot. A little bit different from the PCA result is that, NMF tends to extract separate features instead of mixed ones.

```
1 '''NMF for MNIST'''
2 x_norm = x / x.sum(1).reshape(-1, 1)
3 # Normalize the data
4 Winit = np.random.rand(N, r) + 0.5
5 Hinit = np.random.rand(r, d) + 0.5
6 W, H = nmf(x_norm, Winit, Hinit, tol=1e-3, maxiter=500)
7 print(H.shape)
8 for i in range(r):
9     u = H[i, :]
10    u_img = u.reshape(28, 28)
11    plt.imshow(u_img)
12    plt.title('NMF r=' + str(i + 1))
13    plt.colorbar(shrink=0.83)
14    # plt.show()
15    plt.savefig(f_outpath + 'NMF r= ' + str(i + 1) + '.png', format='png', transparent=True, dpi=300, pad_
16               inches = 0)
17    plt.clf()
18    plt.cla()
```

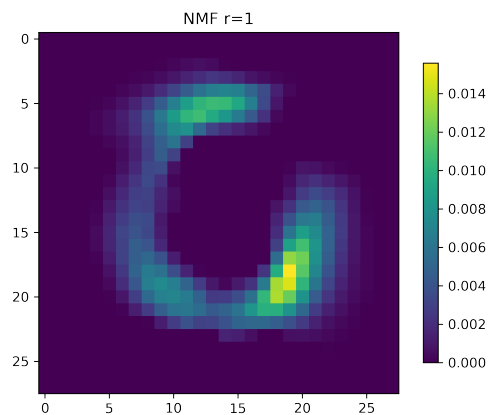


Figure 7: NMF for MNIST:  $r = 1$

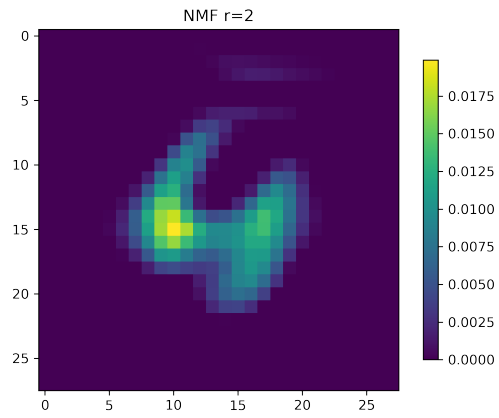


Figure 8: NMF for MNIST:  $r = 2$

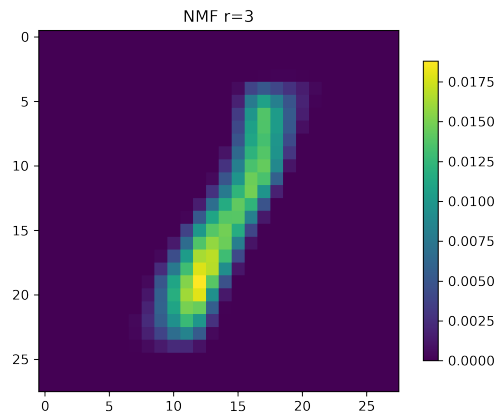


Figure 9: NMF for MNIST:  $r = 3$

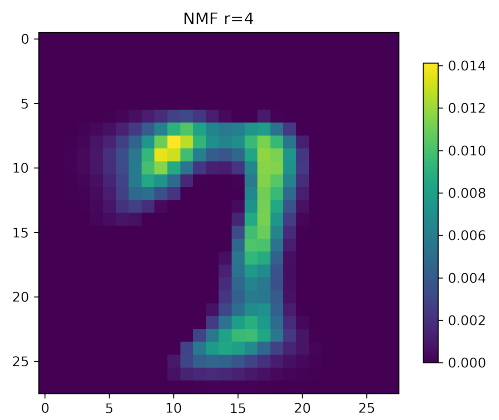


Figure 10: NMF for MNIST:  $r = 4$

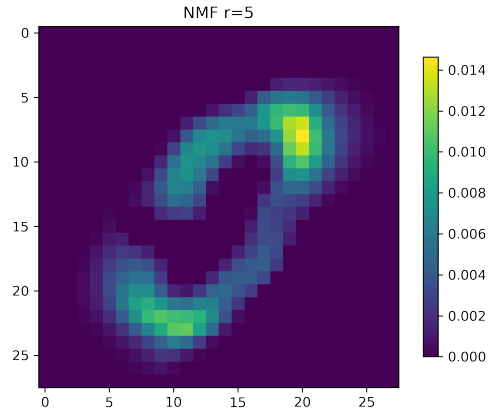


Figure 11: NMF for MNIST:  $r = 5$

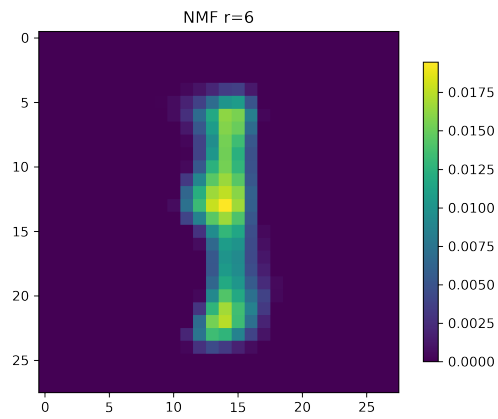


Figure 12: NMF for MNIST:  $r = 6$