

Mathematical Foundations of Machine Learning.

Homework 5

Teresa Morales

November 22nd, 2019

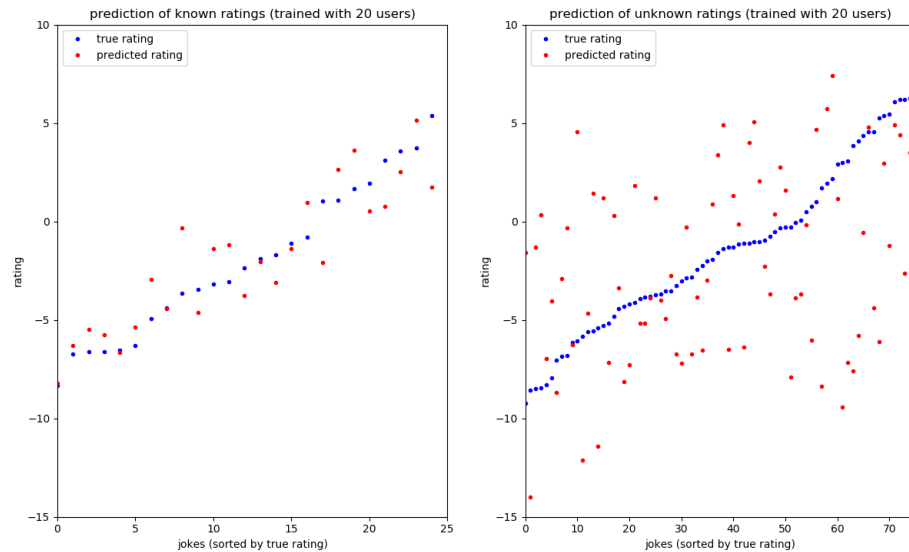
1. Problem 1

(a) Least squares using 20 customers.

The prediction seems to work well with training data but it works much worse when using the testing sample. This might suggest that the sample is too small for good predictions on new data to be made.

Figure 1 shows results graphically:

Figure 1:



The average least squares error for the training data is 1.719.

The average least squares error for the test sample is 5.362

The code used is the following:

```

1 import scipy.io as sio
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from numpy.linalg import norm
5 from mpl_toolkits.mplot3d import Axes3D
6
7
8 # Problem 1 a)
9
10 file_path = '/Users/teresamoraes/Documents/Harris/MFML/
    Homework5/jesterdata.mat'
11 jester_data = sio.loadmat(file_path)
12 X = jester_data['X']
13
14
15 file_path = '/Users/teresamoraes/Documents/Harris/MFML/
    Homework5/newuser.mat'
16 d_new = sio.loadmat(file_path)
17
18 y = d_new['y']
19 true_y = d_new['truey']
20
21 # total number of joke ratings should be m = 100, n = 7200
22 m, n = X.shape
23
24 # train on ratings we know for the new user
25 train_indices = np.squeeze(y != -99)
26 num_train = np.count_nonzero(train_indices)
27
28 # test on ratings we don't know
29 test_indices = np.logical_not(train_indices)
30 num_test = m - num_train
31
32 X_data = X[train_indices, 0:20]
33 y_data = y[train_indices]
34 y_test = true_y[test_indices]
35
36 # solve for weights
37
38
39 def calc_weights(X, y):
40     Xt = np.transpose(X)
41     XtX = np.matmul(Xt, X)
42     Inverse_XtX = np.linalg.inv(XtX)
43     w = np.matmul(np.matmul(Inverse_XtX, Xt), y)
44     return w
45
46
47 w_hat = calc_weights(X_data, y_data)
48
49 # compute predictions
50 X_test = X[test_indices, 0:20]
51 y_hat_train = np.matmul(X_data, w_hat) # Prediction for
    training data
52 y_hat_test = np.matmul(X_test, w_hat) # Prediction for
    test data
53

```

```

54 # measure performance on training jokes
55 error_train = np.subtract(y_data, y_hat_train)
56 error_train_squares = np.square(error_train)
57 avgerr_train = np.sqrt(np.mean(error_train_squares))
58
59 # display results
60
61 ax1 = plt.subplot(121)
62 sorted_indices = np.argsort(np.squeeze(y_data))
63 ax1.plot(range(num_train), y_data[sorted_indices], 'b.',
64         range(num_train), y_hat_train[sorted_indices], 'r.'
65         .')
66 ax1.set_title('prediction of known ratings (trained with
67              20 users)')
68 ax1.set_xlabel('jokes (sorted by true rating)')
69 ax1.set_ylabel('rating')
70 ax1.legend(['true rating', 'predicted rating'], loc='upper
71           left')
72 ax1.axis([0, num_train, -15, 10])
73 print("Average l_2 error (train):", avgerr_train)
74
75 # measure performance on unrated jokes
76 error_test = np.subtract(y_test, y_hat_test)
77 error_test_squares = np.square(error_test)
78 avgerr_test = np.sqrt(np.mean(error_test_squares))
79
80 # display results
81 ax2 = plt.subplot(122)
82 sorted_indices = np.argsort(np.squeeze(y_test))
83 ax2.plot(range(num_test), y_test[sorted_indices], 'b.',
84         range(num_test), y_hat_test[sorted_indices], 'r.'
85         .')
86 ax2.set_title('prediction of unknown ratings (trained with
87              20 users)')
88 ax2.set_xlabel('jokes (sorted by true rating)')
89 ax2.set_ylabel('rating')
90 ax2.legend(['true rating', 'predicted rating'], loc='upper
91           left')
92 ax2.axis([0, num_test, -15, 10])
93 print("Average l_2 (test):", avgerr_test)
94 plt.show()

```

- (b) Using the entire features matrix.

With the new dimensions, there are infinitely many solutions. We can find the least squares solution, which will have the smallest norm.

In order to find the least squares solutions we need to redefine how we estimate the weights using the new dimensions. In particular, before we had:

$$\hat{w} = (X^T X)^{-1} X^T y$$

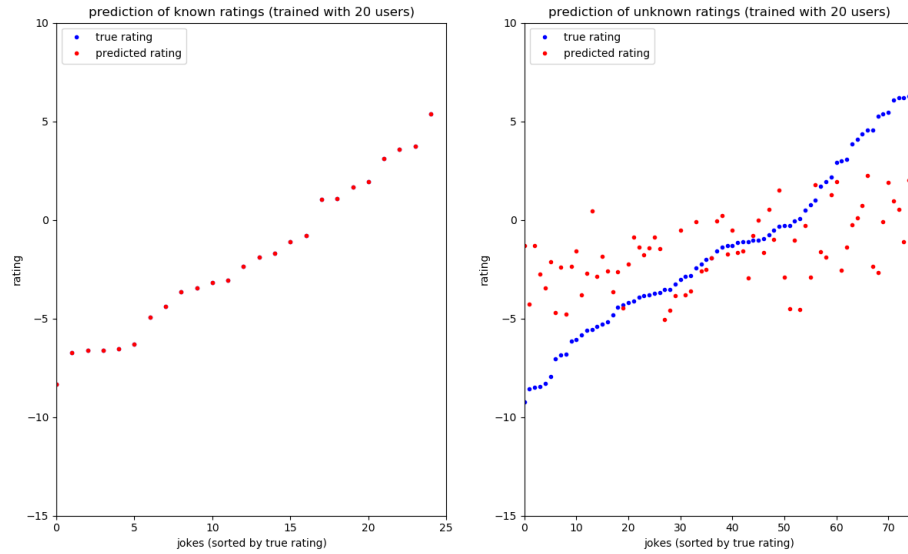
and now we have

$$\hat{w} = X^T (X X^T)^{-1} y$$

With this additional information, errors get smaller. In particular, the mean error for the training sample is 0 and the error for the testing sample is 3.494.

Figure 2 shows results graphically:

Figure 2:



The code used is the following:

```

1 def calc_weights_under(X, y):
2     Xt = np.transpose(X)
3     XXt = np.matmul(X, Xt)
4     Inverse_XXt = np.linalg.inv(XXt)
5     w = np.matmul(np.matmul(Xt, Inverse_XXt), y)
6     return w
7
8
9 X_data = X[train_indices, :]
10 w_hat_under = calc_weights_under(X_data, y_data)
11
12 # compute predictions
13
14 X_test = X[test_indices, :]
15 y_hat_train = np.matmul(X_data, w_hat_under) # Prediction
16 y_hat_test = np.matmul(X_test, w_hat_under)  # Prediction
17
18 # measure performance on training jokes
19 error_train = np.subtract(y_data, y_hat_train)

```

```

20 error_train_squares = np.square(error_train)
21 avgerr_train = np.sqrt(np.mean(error_train_squares))
22
23 # display results
24
25 ax1 = plt.subplot(121)
26 sorted_indices = np.argsort(np.squeeze(y_data))
27 ax1.plot(range(num_train), y_data[sorted_indices], 'b.',
28         range(num_train), y_hat_train[sorted_indices], 'r
29         .')
29 ax1.set_title('prediction of known ratings (trained with
30               20 users)')
30 ax1.set_xlabel('jokes (sorted by true rating)')
31 ax1.set_ylabel('rating')
32 ax1.legend(['true rating', 'predicted rating'], loc='upper
33           left')
33 ax1.axis([0, num_train, -15, 10])
34 print("Average l_2 error (train):", avgerr_train)
35
36 # measure performance on unrated jokes
37 error_test = np.subtract(y_test, y_hat_test)
38 error_test_squares = np.square(error_test)
39 avgerr_test = np.sqrt(np.mean(error_test_squares))
40
41 # display results
42 ax2 = plt.subplot(122)
43 sorted_indices = np.argsort(np.squeeze(y_test))
44 ax2.plot(range(num_test), y_test[sorted_indices], 'b.',
45         range(num_test), y_hat_test[sorted_indices], 'r.'
46         )
46 ax2.set_title('prediction of unknown ratings (trained with
47               20 users)')
47 ax2.set_xlabel('jokes (sorted by true rating)')
48 ax2.set_ylabel('rating')
49 ax2.legend(['true rating', 'predicted rating'], loc='upper
50           left')
50 ax2.axis([0, num_test, -15, 10])
51 print("Average l_2 (test):", avgerr_test)
52 plt.show()

```

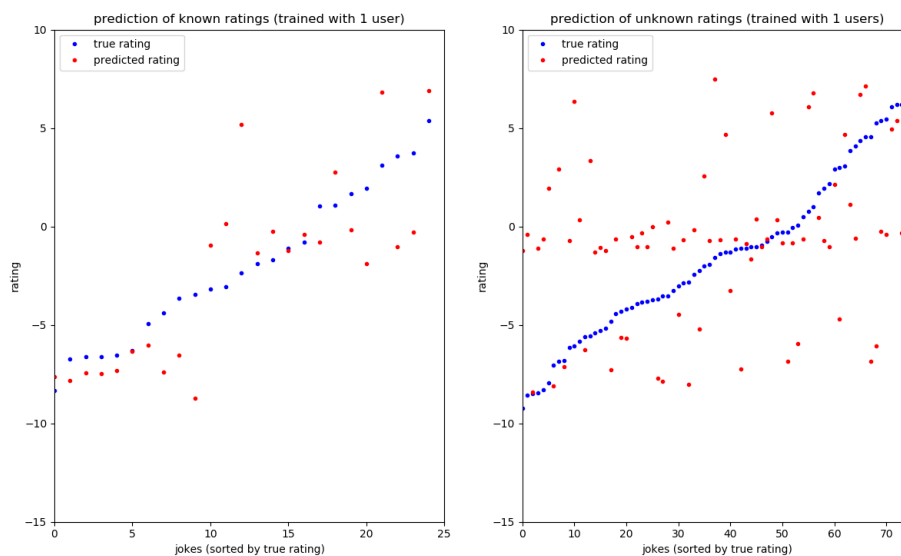
(c) One and two users

Suggested method: estimating the similarity between users using the cosine of the features vectors for the observed ratings.

For one user the mean error is 5.99 for the training sample and 5.82 for the testing samples.

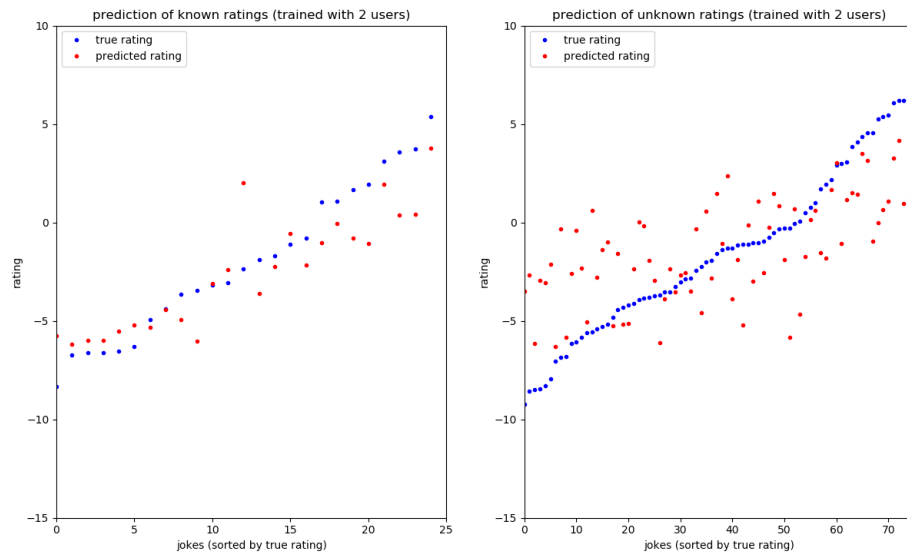
Results are shown in the following graph (figure 3):

Figure 3:



Adding a second user has a huge impact in terms of error. The mean error for the training sample is 1.892 and for the testing sample 3.195. Results are shown in figure 4. For one user the code is the following:

Figure 4:



```

1 def similarity(v1, v2):
2     return np.dot(v1, v2)/(np.linalg.norm(v1)*np.linalg.
3         norm(v2))
4
5 sim_list = [similarity(X_data[:, i], y_data) for i in
6     range(n)]
7
8 cust_1 = sim_list.index(max(sim_list))
9
10 # Customer with index 588 would be the one closest to the
11 # new costumer in ratings of training data
12 y_hat_train = X_data[:, cust_1]
13 y_hat_test = X_test[:, cust_1]
14
15 # measure performance on training jokes
16 error_train = np.subtract(y_data, y_hat_train)
17 error_train_squares = np.square(error_train)
18 avgerr_train = np.sqrt(np.mean(error_train_squares))
19
20 # display results
21 ax1 = plt.subplot(121)
22 sorted_indices = np.argsort(np.squeeze(y_data))
23 ax1.plot(range(num_train), y_data[sorted_indices], 'b.',
24     range(num_train), y_hat_train[sorted_indices], 'r
25     .')
26 ax1.set_title('prediction of known ratings (trained with 1
27     user)')

```

```

24 ax1.set_xlabel('jokes (sorted by true rating)')
25 ax1.set_ylabel('rating')
26 ax1.legend(['true rating', 'predicted rating'], loc='upper
    left')
27 ax1.axis([0, num_train, -15, 10])
28 print("Average l_2 error (train):", avgerr_train)
29
30 # measure performance on unrated jokes
31 error_test = np.subtract(y_test, y_hat_test)
32 error_test_squares = np.square(error_test)
33 avgerr_test = np.sqrt(np.mean(error_test_squares))
34
35 # display results
36 ax2 = plt.subplot(122)
37 sorted_indices = np.argsort(np.squeeze(y_test))
38 ax2.plot(range(num_test), y_test[sorted_indices], 'b.',
39         range(num_test), y_hat_test[sorted_indices], 'r.'
    )
40 ax2.set_title('prediction of unknown ratings (trained with
    1 users)')
41 ax2.set_xlabel('jokes (sorted by true rating)')
42 ax2.set_ylabel('rating')
43 ax2.legend(['true rating', 'predicted rating'], loc='upper
    left')
44 ax2.axis([0, num_test, -15, 10])
45 print("Average l_2 (test):", avgerr_test)
46 plt.show()

```

For two users the code is the following:

```

1 sim_list.remove(max(sim_list))
2
3 cust_2 = sim_list.index(max(sim_list))+1
4
5 X_train_2 = X_data[:, [cust_1, cust_2]]
6 X_test_2 = X_test[:, [cust_1, cust_2]]
7
8 w_hat = calc_weights(X_train_2, y_data)
9
10 # compute predictions
11
12 y_hat_train = np.matmul(X_train_2, w_hat) # Prediction
    for training data
13 y_hat_test = np.matmul(X_test_2, w_hat) # Prediction for
    test data
14
15 # measure performance on training jokes
16 error_train = np.subtract(y_data, y_hat_train)
17 error_train_squares = np.square(error_train)
18 avgerr_train = np.sqrt(np.mean(error_train_squares))
19
20 # display results
21 ax1 = plt.subplot(121)
22 sorted_indices = np.argsort(np.squeeze(y_data))
23 ax1.plot(range(num_train), y_data[sorted_indices], 'b.',
24         range(num_train), y_hat_train[sorted_indices], 'r
    .')
25 ax1.set_title('prediction of known ratings (trained with 2

```



```

        users)')
26 ax1.set_xlabel('jokes (sorted by true rating)')
27 ax1.set_ylabel('rating')
28 ax1.legend(['true rating', 'predicted rating'], loc='upper
    left')
29 ax1.axis([0, num_train, -15, 10])
30
31 print("Average l_2 error (train):", avgerr_train)
32
33 # measure performance on unrated jokes
34 error_test = np.subtract(y_test, y_hat_test)
35 error_test_squares = np.square(error_test)
36 avgerr_test = np.sqrt(np.mean(error_test_squares))
37
38 # display results
39 ax2 = plt.subplot(122)
40 sorted_indices = np.argsort(np.squeeze(y_test))
41 ax2.plot(range(num_test), y_test[sorted_indices], 'b.',
42         range(num_test), y_hat_test[sorted_indices], 'r.'
43         )
44 ax2.set_title('prediction of unknown ratings (trained with
    2 users)')
45 ax2.set_xlabel('jokes (sorted by true rating)')
46 ax2.set_ylabel('rating')
47 ax2.legend(['true rating', 'predicted rating'], loc='upper
    left')
48 ax2.axis([0, num_test, -15, 10])
49 print("Average l_2 (test):", avgerr_test)
50 plt.show()

```

(d) The rank of X is 100.

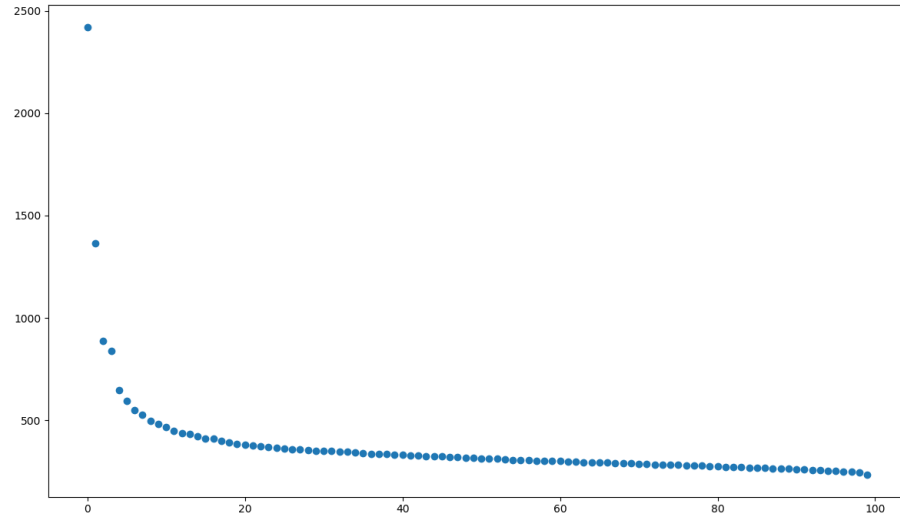
Looking at the plot it seems that only 4 or 5 dimensions seem important. This means that having that number of users, if selected to be as similar to the costumer we are predicting as possible, should give enough information to perform a good prediction. In mathematical terms, this means that we can use a 4 dimension subspace that approximates this data.

```

1 u, s, vh = np.linalg.svd(X, full_matrices=False)
2 s.shape
3 # The rank of X is 100.
4 plt.scatter(range(len(s)), s)
5 plt.show()

```

Figure 5:



(e) Section e)

The first three principal components X (100×7200) are the first three columns of U . The first three principal components of X transpose (7200×100) are the first three columns of U .

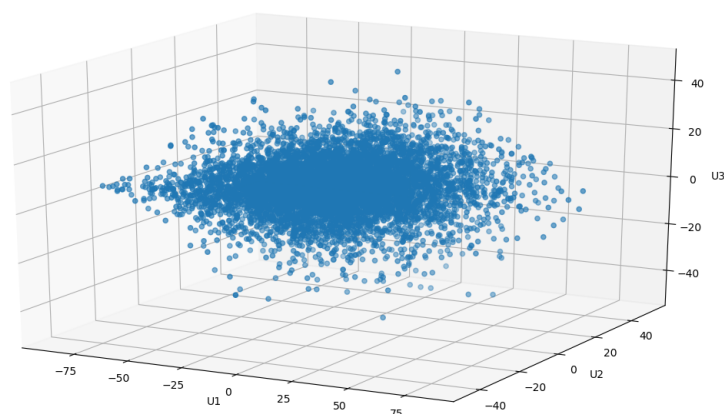
In order to find the projection of each point (each costumer) onto the principal components we would do the following:

```
1 X_p1 = np.matmul(np.diag(s[0:3]), vh[0:3, :])
2
3
4 fig = plt.figure()
5 ax = fig.add_subplot(111, projection='3d')
6 ax.scatter(X_p1[0, :], X_p1[1, :], X_p1[2, :])
7 ax.set_xlabel('U1', fontsize=10)
8 ax.set_ylabel('U2', fontsize=10)
9 ax.set_zlabel('U3', fontsize=10)
10 ax.legend()
11 plt.show()
```

The resulting graph can be seen in figure 6:

As we can see, there is little structure in the data, which means that the three principal components give a good amount of information about costumers' tastes. However, we can still observe some higher variance on one direction which might indicate that a fourth principal

Figure 6:



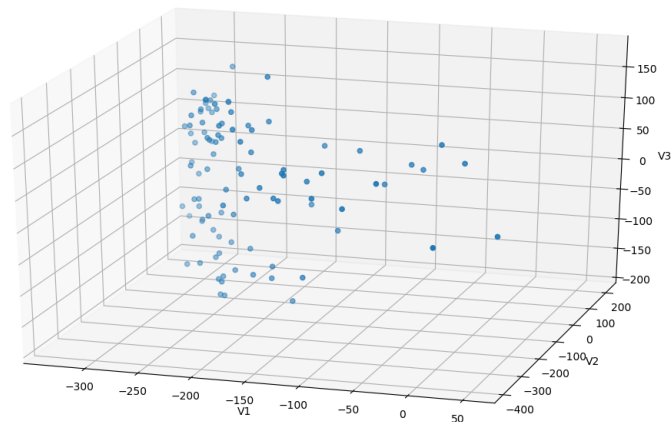
component could help better predict ratings.

Similarly, we can project each joke onto the corresponding three principal components as follows:

```
1 U_t = np.transpose(u)
2 X_p2 = np.matmul(np.diag(s[0:3]), U_t[0:3, :])
3
4 fig = plt.figure()
5 ax = fig.add_subplot(111, projection='3d')
6 ax.scatter(X_p2[0, :], X_p2[1, :], X_p2[2, :])
7 ax.set_xlabel('V1', fontsize=10)
8 ax.set_ylabel('V2', fontsize=10)
9 ax.set_zlabel('V3', fontsize=10)
10 ax.legend()
11 plt.show()
```

The corresponding graph is in figure 7:

Figure 7:



As before, it seems that a fourth principal component could be helpful although probably enough information is captured by the first three principal components to make good predictions.

(f) Power iterations

As we can see the method gives the same result as the `svd` function (except for the sign which is interchangeable across U and V).

```

1 def power_iteration(A, epsilon):
2     u_0 = np.random.rand(A.shape[1])
3     u_i = np.dot(A, u_0)/np.linalg.norm(np.dot(A, u_0))
4     while(np.linalg.norm(u_0-u_i) > epsilon):
5         u_0 = u_i
6         u_i = np.matmul(A, u_0)/np.linalg.norm(np.matmul(A
7             , u_0))
8     return u_i
9
10 XtX = np.matmul(np.transpose(X), X)
11
12 XXt = np.matmul(X, np.transpose(X))
13
14 u_1 = power_iteration(XtX, 0.0000000001)
15
16 similarity(u_1, vh[0, :])
17

```

```

18 # -0.9999999999999998
19
20 u_2 = power_iteration(XXt, 0.0000000001)
21 similarity(u_2, u[:, 0])
22
23 # -1

```

- (g) Our initial guess can't be completely orthogonal to the final vector. This means, that an initial guess that would not work is the second eigenvector (second column of U to find U1 and second column of V when you are trying to find V1).

2. Problem 2

- (a) Truncated SVD solution Results: The average error is 0.0692

Code used:

```

1 import random
2 import scipy.io as sio
3 import numpy as np
4 import matplotlib.pyplot as plt
5
6 file_path = '/Users/teresamorales/Documents/Harris/MFML/
  Homework2/face_emotion_data.mat'
7 face_data = sio.loadmat(file_path)
8 X = face_data['X']
9 y = face_data['y']
10
11
12 # a)
13
14 # Inverse
15
16 def calc_inverse(X, k):
17     n, p = X.shape
18     u, s, vh = np.linalg.svd(X, full_matrices=True)
19     v = np.transpose(vh)
20     u_t = np.transpose(u)
21     sigma_inv_0 = np.array([[0.0 for i in range(n)] for j
  in range(p)])
22     for i in range(k):
23         sigma_inv_0[i, i] = 1/s[i]
24     return np.matmul(np.matmul(v, sigma_inv_0), u_t)
25
26
27 # Randomize sample partition
28
29 def random_partition(sample_size, n_chunks, chunk_size):
30     sample_size_index = list(range(sample_size))
31     random.Random(222).shuffle(sample_size_index) #
  randomizes selection setting seed 222
32     return [sample_size_index[round(chunk_size * i):round(
  chunk_size * (i + 1))]] for i in range(n_chunks)]
33
34
35 def sign(num):

```

```

36     return -1 if num < 0 else 1
37
38
39 # Predict error for sample reserved (both for choosing k
    and for hold-out sample)
40
41 def calc_error(reserved_calc, reserved_exclude,
    features_matrix, labels_vector, n_chunks, chunk_size,
    k):
42     X_reserved_calc = features_matrix[sliced_index[
    reserved_calc]]
43     y_reserved_calc = labels_vector[sliced_index[
    reserved_calc]]
44     if reserved_calc < reserved_exclude:
45         X_training = features_matrix[sum(
46             sliced_index[0:reserved_calc] + sliced_index[
    reserved_calc+1:reserved_exclude] + sliced_index[
    reserved_exclude+1:n_chunks], [])]
47         y_training = labels_vector[sum(sliced_index[0:
    reserved_calc] + sliced_index[reserved_calc +
48
49                                     1:reserved_exclude] +
    sliced_index[reserved_exclude+1:n_chunks], [])]
50     elif reserved_calc > reserved_exclude:
51         X_training = features_matrix[sum(
52             sliced_index[0:reserved_exclude] +
    sliced_index[reserved_exclude+1:reserved_calc] +
    sliced_index[reserved_calc+1:n_chunks], [])]
53         y_training = labels_vector[sum(
54             sliced_index[0:reserved_exclude] +
    sliced_index[reserved_exclude+1:reserved_calc] +
    sliced_index[reserved_calc+1:n_chunks], [])]
55     svd_inv = calc_inverse(X_training, k)
56     w_training = np.matmul(svd_inv, y_training)
57     y_hat = np.matmul(X_reserved_calc, w_training)
58     y_tilda = [sign(i) for i in y_hat]
59     return np.sum([y_tilda[i] != y_reserved_calc[i] for i
    in range(int(chunk_size))])/chunk_size
60
61 # Choose best k
62 def choose_k(reserved_k, reserved_test, features_matrix,
    labels_vector, n_chunks, chunk_size):
63     error_rates_perk = [calc_error(reserved_k,
    reserved_test, features_matrix,
64                                     labels_vector, n_chunks
    , chunk_size, i) for i in range(9)]
65     min_error_k = error_rates_perk.index(min(
    error_rates_perk))
66     return min_error_k
67
68
69 # Run program repeating the process 56 times
70 sample_size = 128
71 number_of_chunks = 8
72 chunk_size = sample_size/number_of_chunks
73

```

```

74 sliced_index = random_partition(sample_size,
    number_of_chunks, chunk_size)
75
76 error_test = []
77
78 for i in range(number_of_chunks):
79     for j in range(number_of_chunks):
80         if j != i:
81             min_k = choose_k(i, j, X, y, number_of_chunks,
                chunk_size)
82             error_test.append(calc_error(j, i, X, y,
                number_of_chunks, chunk_size, min_k))
83
84 average_error = np.mean(error_test)

```

(b) Ridge regression

Results: The average error is 0.0335

Code used:

```

1 # Problem 2b)
2
3 # We redefine the inverse:
4
5 def calc_inverse(X, k):
6     n, p = X.shape
7     u, s, vh = np.linalg.svd(X, full_matrices=True)
8     v = np.transpose(vh)
9     u_t = np.transpose(u)
10    sigma_inv_0 = np.array([[0.0 for i in range(n)] for j
        in range(p)])
11    for i in range(9):
12        sigma_inv_0[i, i] = s[i]/(s[i]**2+k)
13    return np.matmul(np.matmul(v, sigma_inv_0), u_t)
14
15 # We redefine the calculation of optimal lambda (k)
16
17
18 def choose_k(reserved_k, reserved_test, features_matrix,
    labels_vector, n_chunks, chunk_size, lambda_vals):
19     error_rates_perk = [calc_error(reserved_k,
        reserved_test, features_matrix,
        labels_vector, n_chunks
        , chunk_size, i) for i in lambda_vals]
20     min_error_k = lambda_vals[error_rates_perk.index(min(
        error_rates_perk))]
21     return min_error_k
22
23
24
25 lambda_vals = np.array([0, 0.5, 1, 2, 4, 8, 16])
26
27 error_test_ridge = []
28
29 for i in range(number_of_chunks):
30     for j in range(number_of_chunks):
31         if j != i:
32             min_k = choose_k(i, j, X, y, number_of_chunks,
                chunk_size, lambda_vals)

```

```

33         error_test_ridge.append(calc_error(j, i, X, y,
34                                         number_of_chunks, chunk_size, min_k))
35 average_error = np.mean(error_test_ridge)

```

- (c) Theoretically, results should not change since the additional singular values will be 0. However when we perform the operations error rates might change due to rounding errors. Actually when calculating inverses, small rounding errors for zero values can have a huge impact although we are trying to reduce this impact using truncates and ridge regression.

When we repeat parts a and b using these new features we obtain an average error rate of 0.09154 for truncated SVD and 0.04576 for ridge regression.

The increase in error rate due to this noise seems to be smaller for ridge regression than for truncated SVD.

The code used is the following:

```

1  # Problem 2 c)
2
3  New_X = np.hstack((X, np.matmul(X, np.random.randn(9, 3)))
4                    )
5
6
7  # Part a'
8
9
10 def calc_inverse(X, k):
11     n, p = X.shape
12     u, s, vh = np.linalg.svd(X, full_matrices=True)
13     v = np.transpose(vh)
14     u_t = np.transpose(u)
15     sigma_inv_0 = np.array([[0.0 for i in range(n)] for j
16                             in range(p)])
17     for i in range(k):
18         sigma_inv_0[i, i] = 1/s[i]
19     return np.matmul(np.matmul(v, sigma_inv_0), u_t)
20
21 # Choose best k
22
23 def choose_k(reserved_k, reserved_test, features_matrix,
24             labels_vector, n_chunks, chunk_size):
25     error_rates_perk = [calc_error(reserved_k,
26                                     reserved_test, features_matrix,
27                                     labels_vector, n_chunks,
28                                     chunk_size, i) for i in range(12)] # Should we
29                                     # change this to 12 for the experiment or leave 9 since
30                                     # there is rounding error?
31     min_error_k = error_rates_perk.index(min(
32         error_rates_perk))
33     return min_error_k

```



```

29
30 error_test = []
31
32 for i in range(number_of_chunks):
33     for j in range(number_of_chunks):
34         if j != i:
35             min_k = choose_k(i, j, New_X, y,
36                             number_of_chunks, chunk_size)
37             error_test.append(calc_error(j, i, New_X, y,
38                                         number_of_chunks, chunk_size, min_k))
39
40 average_error = np.mean(error_test)
41
42 # Part b'
43 def calc_inverse(X, k):
44     n, p = X.shape
45     u, s, vh = np.linalg.svd(X, full_matrices=True)
46     v = np.transpose(vh)
47     u_t = np.transpose(u)
48     sigma_inv_0 = np.array([[0.0 for i in range(n)] for j
49                             in range(p)])
50     for i in range(12):
51         sigma_inv_0[i, i] = s[i]/(s[i]**2+k)
52     return np.matmul(np.matmul(v, sigma_inv_0), u_t)
53 # We redefine the calculation of optimal lambda (k)
54
55 def choose_k(reserved_k, reserved_test, features_matrix,
56             labels_vector, n_chunks, chunk_size, lambda_vals):
57     error_rates_perk = [calc_error(reserved_k,
58                                     reserved_test, features_matrix,
59                                     labels_vector, n_chunks,
60                                     chunk_size, i) for i in lambda_vals]
61     min_error_k = lambda_vals[error_rates_perk.index(min(
62         error_rates_perk))]
63     return min_error_k
64
65 error_test_ridge = []
66
67 for i in range(number_of_chunks):
68     for j in range(number_of_chunks):
69         if j != i:
70             min_k = choose_k(i, j, New_X, y,
71                             number_of_chunks, chunk_size, lambda_vals)
72             error_test_ridge.append(calc_error(j, i, New_X,
73                                                 y, number_of_chunks, chunk_size, min_k))
74
75 average_error = np.mean(error_test_ridge)

```

3. Problem 3

- (a) The code is the following. Results are shown in the graph.

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 import pandas as pd
4
5 # deblurring
6
7 n = 500
8 k = 30
9 sigma = 0.01
10
11 # generate random piecewise constant signal
12 w = np.zeros((n, 1))
13 w[0] = np.random.standard_normal()
14 for i in range(1, n):
15     if np.random.rand(1) < 0.95:
16         w[i] = w[i-1]
17     else:
18         w[i] = np.random.standard_normal()
19
20
21 # generate k-point averaging function
22 h = np.ones(k) / k
23
24 # make a matrix for blurring
25 m = n + k - 1
26 X = np.zeros((m, m))
27 for i in range(m):
28     if i < k:
29         X[i, :i+1] = h[:i+1]
30     else:
31         X[i, i - k: i] = h
32
33 X = X[:, 0:n]
34
35 # blurred signal + noise
36 y = np.dot(X, w) + sigma*np.random.standard_normal(size=(m
37 , 1))
38
39 # plot
40 f, (ax1, ax2) = plt.subplots(1, 2)
41 ax1.set_title('signal')
42 ax1.plot(w, 'b')
43 ax2.set_title('blurred and noisy version')
44 ax2.plot(y[0:n])
45
46 plt.show()
47
48 # Problem 3 a)
49 # (1 )Least squares
50
51
52 def calc_weights(X, y):
53     Xt = np.transpose(X)
54     XtX = np.matmul(Xt, X)
55     Inverse_XtX = np.linalg.inv(XtX)
56     w = np.matmul(np.matmul(Inverse_XtX, Xt), y)

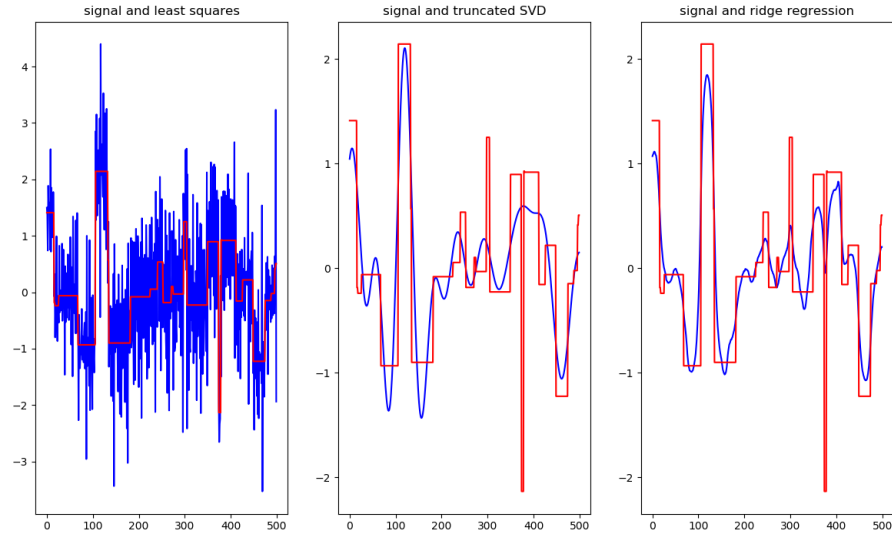
```

```

57     return w
58
59
60 w_hat = calc_weights(X, y)
61
62
63 # (2 )Truncated SVD setting k to 20
64
65 def calc_inverse_trunc(X, k):
66     n, p = X.shape
67     u, s, vh = np.linalg.svd(X, full_matrices=True)
68     v = np.transpose(vh)
69     u_t = np.transpose(u)
70     sigma_inv_0 = np.array([[0.0 for i in range(n)] for j
71                             in range(p)])
72     for i in range(k):
73         sigma_inv_0[i, i] = 1/s[i]
74     return np.matmul(np.matmul(v, sigma_inv_0), u_t)
75
76 SVD_trunc = calc_inverse_trunc(X, 20)
77
78 w_trunc = np.matmul(SVD_trunc, y)
79
80
81 # (3 ) Ridge regression setting lambda to 0.1
82
83 def calc_inverse_ridge(X, k):
84     n, p = X.shape
85     u, s, vh = np.linalg.svd(X, full_matrices=True)
86     v = np.transpose(vh)
87     u_t = np.transpose(u)
88     sigma_inv_0 = np.array([[0.0 for i in range(n)] for j
89                             in range(p)])
90     for i in range(500):
91         sigma_inv_0[i, i] = s[i]/(s[i]**2+k)
92     return np.matmul(np.matmul(v, sigma_inv_0), u_t)
93
94 SVD_ridge = calc_inverse_ridge(X, 0.1)
95
96 w_ridge = np.matmul(SVD_ridge, y)
97
98 f, (ax1, ax2, ax3) = plt.subplots(1, 3)
99 ax1.set_title('signal and least squares')
100 ax1.plot(w_hat, 'b')
101 ax1.plot(w, 'r')
102 ax2.set_title('signal and truncated SVD')
103 ax2.plot(w_trunc, 'b')
104 ax2.plot(w, 'r')
105 ax3.set_title('signal and ridge regression')
106 ax3.plot(w_ridge, 'b')
107 ax3.plot(w, 'r')
108
109 plt.show()

```

Figure 8:



- (b) Using the suggested values for λ and k , the results we obtain can be observed in figure 9.

As we can see in the left graph for a constant σ , the optimal truncating k gets smaller as the blurring k increases. However for extreme values of σ , the relation is constant (the optimal k does not change when σ is as small as 0.01 or as big as 10).

Similarly, in that same graph we observe that, for a constant blurring k , the optimal truncating k gets smaller as the value of σ increases (see shifts in the lines as σ changes).

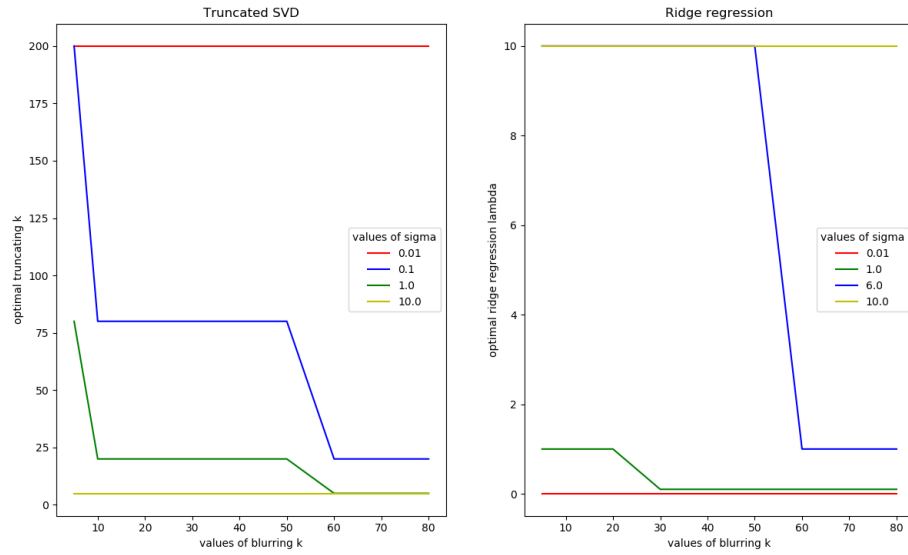
Looking at the right graph we observe that:

- Similarly to what we saw for truncated SVD, the optimal λ gets smaller as the blurring k increases (for a constant σ). As before, for extreme values of σ the relation is constant.
- However, contrary to what we saw before, for a constant blurring k , the optimal λ increases as σ increases.

The code used to produce that graph is the following.

```
1 def calc_min_error_trunc(X, y):
2     list_trunc = [5, 20, 80, 200]
3     error_sq = []
4     for k in list_trunc:
5         SVD_trunc = calc_inverse_trunc(X, k)
6         w_trunc = np.matmul(SVD_trunc, y)
7         error_sq.append(sum((w_trunc-w)**2))
```

Figure 9:



```

8     k_min_error = list_trunc[error_sq.index(min(error_sq))
9     ]
10     return k_min_error
11
12 def calc_min_error_ridge(X, y):
13     list_ridge = [0.01, 0.1, 1, 10]
14     error_sq = []
15     for k in list_ridge:
16         SVD_ridge = calc_inverse_ridge(X, k)
17         w_ridge = np.matmul(SVD_ridge, y)
18         error_sq.append(sum((w_ridge-w)**2))
19     k_min_error = list_ridge[error_sq.index(min(error_sq))
20     ]
21     return k_min_error
22
23 k_list = [5, 10, 20, 30, 40, 50, 60, 80]
24 sigma = [0.000001, 0.0001, 0.01, 0.1, 1, 2, 6, 10]
25
26 opt_trunc = {s: {k: 0 for k in k_list} for s in sigma}
27 opt_ridge = {s: {k: 0 for k in k_list} for s in sigma}
28
29 for s in sigma:
30     for k in k_list:
31         h = np.ones(k) / k
32         m = n + k - 1
33         X = np.zeros((m, m))

```

```

34         for i in range(m):
35             if i < k:
36                 X[i, :i+1] = h[:i+1]
37             else:
38                 X[i, i - k: i] = h
39             X = X[:, 0:n]
40             y = np.dot(X, w) + s*np.random.standard_normal(
size=(m, 1))
41             opt_trunc[s][k] = calc_min_error_trunc(X, y)
42             opt_ridge[s][k] = calc_min_error_ridge(X, y)
43
44 trunc_data = pd.DataFrame(data=opt_trunc)
45 ridge_data = pd.DataFrame(data=opt_ridge)
46
47 sigma = [0.000001, 0.0001, 0.01, 0.1, 1, 2, 6, 10]
48
49 f, (ax1, ax2) = plt.subplots(1, 2)
50 ax1.set_title('Truncated SVD')
51 ax1.plot(trunc_data[0.01], 'r')
52 ax1.plot(trunc_data[0.1], 'b')
53 ax1.plot(trunc_data[1], 'g')
54 ax1.plot(trunc_data[10], 'y')
55 ax1.legend(title='values of sigma')
56 ax1.set_xlabel('values of blurring k')
57 ax1.set_ylabel('optimal truncating k')
58 ax2.set_title('Ridge regression')
59 ax2.plot(ridge_data[0.01], 'r')
60 ax2.plot(ridge_data[1], 'g')
61 ax2.plot(ridge_data[6], 'b')
62 ax2.plot(ridge_data[10], 'y')
63 ax2.legend(title='values of sigma')
64 ax2.set_xlabel('values of blurring k')
65 ax2.set_ylabel('optimal ridge regression lambda')
66
67 plt.show()

```

4. Problem 4

(a) Let

$$x = \begin{bmatrix} 1 \\ 1 \\ \dots \\ 1 \end{bmatrix}$$

$$A^T x = \begin{bmatrix} a_1^T * 1 \\ a_2^T * 1 \\ \dots \\ a_n^T * 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ \dots \\ 1 \end{bmatrix}$$

because each column of A sums to one. This is to say, $A^T x = x$
 So x is an eigenvector and 1 is an eigenvalue of A^T and A.

(b) By inspection we observe that to keep the columns of G summing to one, the elements of vector u need to also sum to 1.

(c)

$$M = \begin{bmatrix} 1 & 1 & 1 & . & . & 1 \\ 0 & 0 & 0 & . & . & 0 \\ 0 & 0 & 0 & . & . & 0 \\ 0 & 0 & 0 & . & . & 0 \\ 0 & 0 & 0 & . & . & 0 \\ 0 & 0 & 0 & . & . & 0 \end{bmatrix}$$

$$G = \begin{bmatrix} \frac{(n-1)\alpha+1}{\frac{n}{(1-\alpha)}} & \frac{(n-1)\alpha+1}{\frac{n}{(1-\alpha)}} & . & . & . & \frac{(n-1)\alpha+1}{\frac{n}{(1-\alpha)}} \\ . & . & . & . & . & . \\ . & . & . & . & . & . \\ . & . & . & . & . & . \\ \frac{(1-\alpha)}{n} & \frac{(1-\alpha)}{n} & . & . & . & \frac{(1-\alpha)}{n} \end{bmatrix}$$

We calculate the eigenvector π as follows:

$$G\pi = \pi$$

which implies,

$$(\pi_1 + \pi_2 + \dots + \pi_n) \frac{(n-1)\alpha+1}{n} = \pi_1$$

and

$$(\pi_1 + \pi_2 + \dots + \pi_n) \frac{(1-\alpha)}{n} = \pi_i$$

for all $i \neq 1$

Since $(\pi_1 + \pi_2 + \dots + \pi_n) = 1$ the Pagerank of Facebook is:

$$\frac{(n-1)\alpha+1}{n} = \pi_1$$

And the PageRank of the rest $n-1$ Web pages is:

$$\frac{(1-\alpha)}{n} = \pi_i$$

for all $i \neq 1$

We can check that the resulting probability vector sums to 1 as follows:

$$\frac{(n-1)\alpha+1}{n} + (n-1) \frac{(1-\alpha)}{n} = \frac{(1-\alpha+\alpha)(n-1)+1}{n} = \frac{n}{n} = 1$$

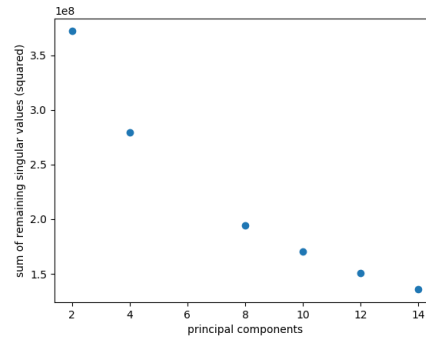
5. Problem 5

(a) Reconstruction accuracy The code used is the following:

```
1 import random
2 import scipy.io as sio
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import pandas as pd
6
7
8 file_path = '/Users/teresamorales/Documents/Harris/MFML/
   Homework5/mnist.mat'
9 number_data = sio.loadmat(file_path)
10 train_data = number_data['train_data']
11 test_data = number_data['test_data']
12 y_train = number_data['train_target']
13 y_test = number_data['test_target']
14
15 u, s, vh = np.linalg.svd(train_data, full_matrices=False)
16
17 list_pc = [2, 4, 8, 10, 12, 14]
18
19 s_squared = [s[i]**2 for i in range(200)]
20
21 acc = {j: sum(s_squared) - sum([s_squared[i] for i in
   range(j)]) for j in list_pc}
22
23 accuracy = {'principal_components': list(acc.keys()),
24            'reconstruction_accuracy': list(acc.values())}
25
26 accuracy_df = pd.DataFrame(data=accuracy)
27
28 fig, ax = plt.subplots()
29 ax.scatter(accuracy_df['principal_components'],
   accuracy_df['reconstruction_accuracy'])
30 ax.set_xlabel("principal components")
31 ax.set_ylabel("sum of remaining singular values (squared)"
   )
32 plt.show()
```

The resulting plot is in figure 10:

Figure 10:



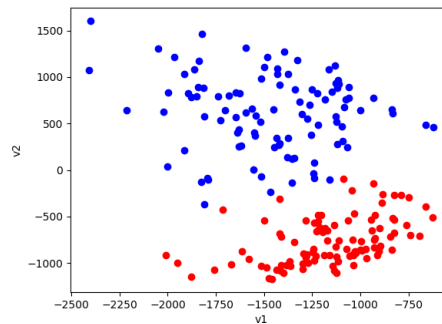
(b) Projection into 2d space

The code used is the following:

```
1 u_t = np.transpose(u)
2
3 X_p = np.matmul(np.diag(s[0:2]), u_t[0:2, :])
4
5 fig, ax = plt.subplots()
6 ax.scatter(X_p[0, 0:100, :], X_p[1, 0:100, :], color='r')
7 ax.scatter(X_p[0, 100:200, :], X_p[1, 100:200, :], color='b')
8 ax.set_xlabel("v1")
9 ax.set_ylabel("v2")
10 plt.show()
```

The resulting plot is the following

Figure 11:



(c) Linear regression

The code used is the following:

```

1
2 def calc_weights(X, y):
3     Xt = np.transpose(X)
4     XtX = np.matmul(Xt, X)
5     Inverse_XtX = np.linalg.inv(XtX)
6     w = np.matmul(np.matmul(Inverse_XtX, Xt), y)
7     return w
8
9
10 col1 = np.ones((200, 1))
11 X_p_constant = np.hstack((np.transpose(X_p), col1))
12
13 w_hat = calc_weights(X_p_constant, np.transpose(y_train))
14
15
16 # measure performance on training sample
17 y_hat_train = np.matmul(X_p_constant, w_hat)
18 error_train = np.subtract(y_train, y_hat_train)
19 error_train_squares = np.square(error_train)
20 avgerr_train = np.sqrt(np.mean(error_train_squares))
21
22 print('Mean squared error for training sample is',
23       avgerr_train)
24
25 # measure performance on test sample
26
27 # project test data into the same subspace as the train
28 # data (using vh from train data )
29 X_test_p = np.matmul(vh[0:2, :], np.transpose(test_data))
30
31 X_test_constant = np.hstack((np.transpose(X_test_p), col1)
32 )
33 y_hat_test = np.matmul(X_test_constant, w_hat) #
34 # Prediction for test data
35
36 error_test = np.subtract(y_test, y_hat_test)
37 error_test_squares = np.square(error_test)
38 avgerr_test = np.sqrt(np.mean(error_test_squares))
39
40 print('Mean squared error for test sample is', avgerr_test
41 )

```

Results: Mean squared error for training sample is 1.3482749026367424

Mean squared error for test sample is 1.3528216855250104

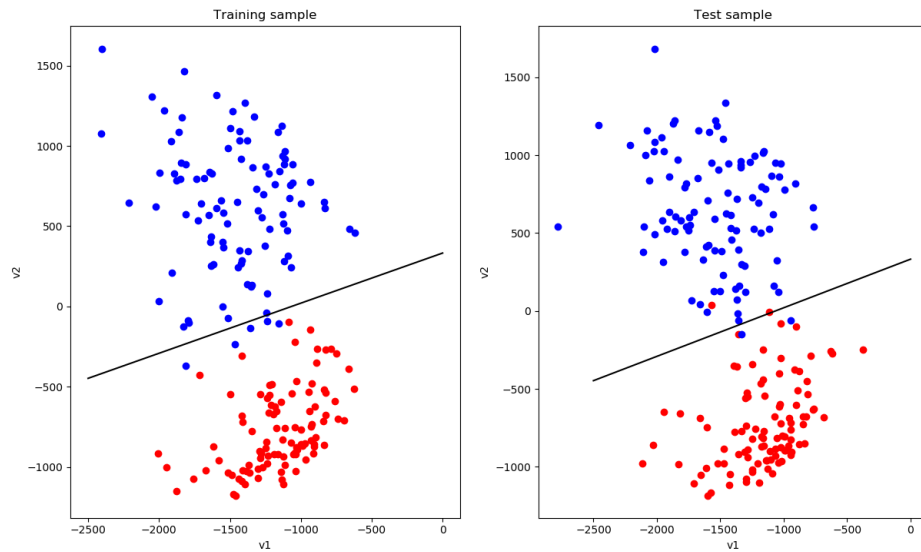
(d) The code and plot would be the following:

```

1 f, (ax1, ax2) = plt.subplots(1, 2)
2 ax1.scatter(X_p[0, 0:100, :], X_p[1, 0:100], color='r')
3 ax1.scatter(X_p[0, 100:200], X_p[1, 100:200], color='b')
4 ax1.plot([-2500, 0], [(-w_hat[0][0]*-2500/w_hat[1][0]) -
5                       w_hat[2]
6                       [0]/w_hat[1][0], -w_hat[2][0]/w_hat
7                       [1][0]], color='black')
8 ax1.set_xlabel("v1")
9 ax1.set_ylabel("v2")

```

Figure 12:



```

8 ax1.set_title('Training sample')
9 ax2.scatter(X_test_p[0, 0:100, :], X_test_p[1, 0:100],
10             color='r')
11 ax2.scatter(X_test_p[0, 100:200], X_test_p[1, 100:200],
12             color='b')
13 ax2.plot([-2500, 0], [(-w_hat[0][0]*-2500/w_hat[1][0]) -
14                       w_hat[2]
15                       [0]/w_hat[1][0], -w_hat[2][0]/w_hat
16                       [1][0]], color='black')
17 ax2.set_xlabel("v1")
18 ax2.set_ylabel("v2")
19 ax2.set_title('Test sample')
20 plt.show()

```

6. Problem 6

- (a) Only two dimensions are necessary since the third singular value is 0 (very close to zero given rounding error)

```

1 import random
2 import scipy.io as sio
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import pandas as pd
6
7

```

```

8 file_path = '/Users/teresamoraless/Documents/Harris/MFML/
  Homework5/pca_3d.mat'
9 p6_data = sio.loadmat(file_path)
10 X = p6_data['point']
11 y = p6_data['target']
12
13 u, s, vh = np.linalg.svd(X, full_matrices=True)

```

- (b) Projection into 2 and 1 dimensional spaces For the projection onto a 2 dimensional space we would use the following code:

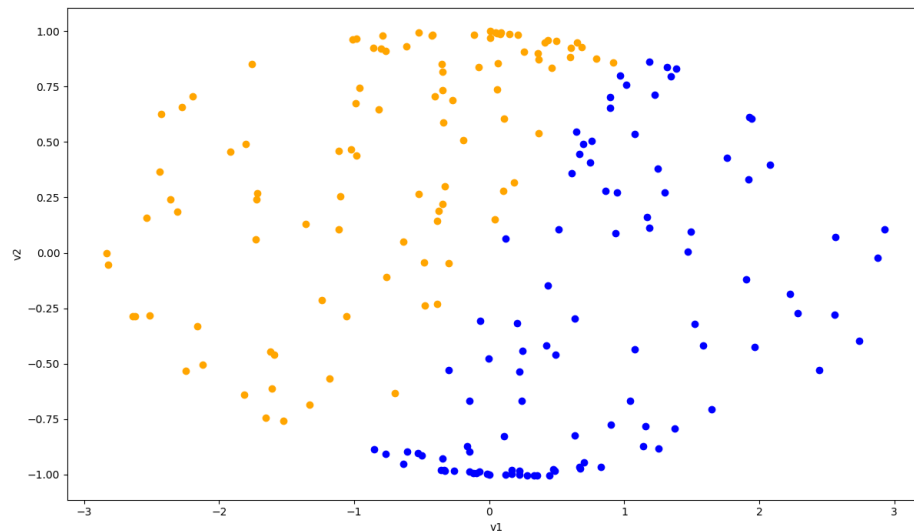
```

1 u_t = np.transpose(u)
2 X_p = np.matmul(np.diag(s[0:2]), u_t[0:2, :])
3
4 #graph
5 value_1 = np.squeeze(y == 1)
6 value_0 = np.logical_not(value_1)
7
8 fig, ax = plt.subplots()
9 ax.scatter(X_p[0, value_1], X_p[1, value_1], color='orange')
10 ax.scatter(X_p[0, value_0], X_p[1, value_0], color='b')
11 ax.set_xlabel("v1")
12 ax.set_ylabel("v2")
13 plt.show()

```

The resulting plot is the following:

Figure 13:

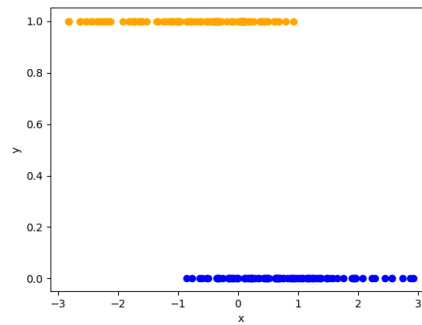


For one dimensions the code would be the following:

```
1 X_p = np.matmul(np.diag(s[0:1]), u_t[0:1, :])
```

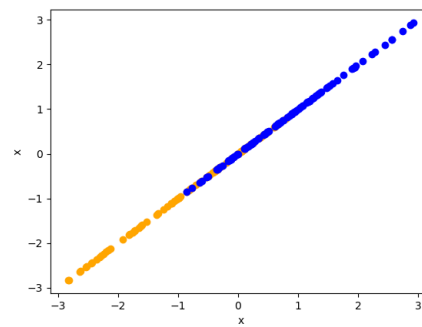
In order to graph this new projection we represent the features in the x axis and the labels in the y axis.

Figure 14:



Another possible way to represent the subspace is plotting the features against their own.

Figure 15:



(c) After whitening the features matrix, the resulting plot is the same as in part b for two dimensions, except for the scale in the axes.

```

1 W_pca = np.matmul(np.diag(1/s[0:2]), vh[0:2, :])
2
3 New_X = np.matmul(X, np.transpose(W_pca))
4
5
6 def check_diag(X_tilda):
7     return np.matmul(np.transpose(X_tilda), X_tilda)
8
9
10 fig, ax = plt.subplots()
11 ax.scatter(New_X[value_1, 0], New_X[value_1, 1], color='
    orange')
12 ax.scatter(New_X[value_0, 0], New_X[value_0, 1], color='b'
    )
13 ax.set_xlabel("v1")
14 ax.set_ylabel("v2")
15 plt.show()

```

Figure 16:

