**Chang'an University**

# Exercise Report

## *"K*-Means Clustering and Principal Component Analysis"*

**May 04, 2022**

**Submitted By:**
Manik Debnath |2021124913
Computer Science and Technology
School of Information Engineering
Chang'an University
Statistical Machine Learning

**Submitted To:**
Prof. Qingsong Song
School of Information Engineering
Chang'an University

# <mark>Machine Learning: Programming Exercise 7</mark>

## *K*-Means Clustering and Principal Component Analysis

## Introduction:

Machine learning algorithms can be broadly classified into two categories - supervised and unsupervised learning. There are other categories also like semi-supervised learning and reinforcement learning. But, most of the algorithms are classified as supervised or unsupervised learning. The difference between them happens because of presence of target variable. In unsupervised learning, there is no target variable. The dataset only has input variables which describe the data. This is called unsupervised learning.

Clustering is the method of dividing a set of abstract objects into groups. Points to Keep in Mind **a set of data objects can be viewed as a single entity**. When performing cluster analysis, we divide the data set into groups based on data similarity, then assign labels to the groups. The various types of clustering are Hierarchical clustering and Partitioning clustering. Partitioning clustering is further subdivided into K-Means clustering and Fuzzy C-Means clustering.

K-Means clustering is the most popular unsupervised learning algorithm. It is used when we have unlabeled data which is data without defined categories or groups. The algorithm follows an easy or simple way to classify a given data set through a certain number of clusters, fixed apriori. K-Means algorithm works iteratively to assign each data point to one of K groups based on the features that are provided. Data points are clustered based on feature similarity.

**K-Means clustering** is used to find intrinsic groups within the unlabeled dataset and draw inferences from them. It is based on centroid-based clustering.

**Principal Component Analysis (PCA)** is a technique that is widely used for applications such as dimensionality reduction, visualization and lossy data compression. In this lab we will focus on the visualization aspect. PCA can be defined as the orthogonal linear transformation of the data to a lower dimensional linear space, known as the principal subspace, such that the greatest variance by any projection of the data comes to lie on the first coordinate (called the first principal component), the second greatest variance on the second coordinate, and so on. Intuitively, PCA finds a meaningful coordinate basis to express our data.

**In this exercise, you will implement the K-means clustering algorithm and apply it to compress an image. In the second part, you will use principal component analysis to find a low-dimensional representation of face images.**

# TABLE OF CONTENTS

# LIST OF FIGURES

# Part 1: K-means Clustering

In this this exercise, we will implement the *K*-means algorithm and use it for image compression. We will first start on an example 2D dataset that will help you gain an intuition of how the *K*-means algorithm works. After that, you will use the *K*-means algorithm for image compression by reducing the number of colors that occur in an image to only those that are most common in that image.

## Part 1.1: Implementing K-means

The K-means algorithm is a method to automatically cluster similar data examples together. Concretely, you are given a training set $\{x^{(1)}, \ldots x^{(m)}\}$ (where $x^{(i)} \in \mathbb{R}^n$), and want to group the data into a few cohesive 'clusters'. The intuition behind *K*-means is an iterative procedure that starts by guessing the initial centroids, and then refines this guess by repeatedly assigning examples to their closest centroids and then recomputing the centroids based on the assignments.

**The *K*-means algorithm is as follows:**

```
% Initialize centroids
centroids = kMeansInitCentroids(X, K);
for iter = 1:iterations
    % Cluster assignment step: Assign each data point to the
    % closest centroid. idx(i) corresponds to c^(i), the index
    % of the centroid assigned to example i
    idx = findClosestCentroids(X, centroids);

    % Move centroid step: Compute means based on centroid
    % assignments
    centroids = computeMeans(X, idx, K);
end
```

*Figure 1: K-means algorithm*

The inner-loop of the algorithm repeatedly carries out two steps:

1. Assigning each training example $x^{(i)}$ to its closest centroid
2. Recomputing the mean of each centroid using the points assigned to it.

The *K*-means algorithm will always converge to some final set of means for the centroids. Note that the converged solution may not always be ideal and depends on the initial setting of the centroids. Therefore, in practice the K-means algorithm is usually run a few times with different random initializations. One way to choose between these different solutions from different random initializations is to choose the one with the lowest cost function value (distortion). You will implement the two phases of the K-means algorithm separately in the next sections.

## Part 1.1.1: Finding closest centroids

In the 'cluster assignment' phase of the K-means algorithm, the algorithm assigns every training example $x^{(i)}$ to its closest centroid, given the current positions of centroids. Specially, for every example $i$ we set

$$c^{(i)} := j \quad \text{that minimizes} \quad ||x^{(i)} - \mu_j||^2,$$

where $c^{(i)}$ is the index of the centroid that is closest to $x^{(i)}$, and $\mu_j$ is the position (value) of the $j$-th centroid. Note that $c^{(i)}$ corresponds to idx (i) in the starter code.

ex7.m

```
ⓘ This file uses Cell Mode. For information, see the rapid code iteration video, the publishing video, or help.        ×
23
24    %% ================= Part 1: Find Closest Centroids =====================
25    %  To help you implement K-Means, we have divided the learning algorithm
26    %  into two functions -- findClosestCentroids and computeCentroids. In this
27    %  part, you shoudl complete the code in the findClosestCentroids function.
28    %
29 -  fprintf('Finding closest centroids.\n\n');
30
31    % Load an example dataset that we will be using
32 -  load('ex7data2.mat');
33
34    % Select an initial set of centroids
35 -  K = 3; % 3 Centroids
36 -  initial_centroids = [3 3; 6 2; 8 5];
37
38    % Find the closest centroids for the examples using the
39    % initial_centroids
40 -  idx = findClosestCentroids(X, initial_centroids);
41
42 -  fprintf('Closest centroids for the first 3 examples: \n')
43 -  fprintf(' %d', idx(1:3));
44 -  fprintf('\n(the closest centroids should be 1, 3, 2 respectively)\n');
45
46 -  fprintf('Program paused. Press enter to continue.\n');
47 -  pause;
48
```
displayData.m  ×  drawLine.m  ×  ex7.m  ×  ex7.mlx  ×  ex7_companion.mlx  ×  ex7_pca.m  ×  featureNormalize.m  ×  findClosestCentroids.m  ×  k ▶

script                              Ln 1      Col 1      OVR

*Figure 2: Find closest centroids (ex7 p1)*

Our task is to complete the code in **findClosestCentroids.m**. This function takes the data matrix X and the locations of all centroids inside centroids and should output a one-dimensional array idx that holds the index (a value in $\{1, \dots, K\}$, where $K$ is total number of centroids) of the closest centroid to every training example. We can implement this using a loop over every training example and every centroid.

Once we have completed the code in **findClosestCentroids.m**, the code below will run our code and we should see the output [1 3 2] corresponding to the centroid assignments for the first 3 examples.

findClosestCentroids.m

```matlab
function idx = findClosestCentroids(X, centroids)
%FINDCLOSESTCENTROIDS computes the centroid memberships for every example
%   idx = FINDCLOSESTCENTROIDS (X, centroids) returns the closest centroids
%   in idx for a dataset X where each row is a single example. idx = m x 1
%   vector of centroid assignments (i.e. each entry in range [1..K])
%

% Set K
K = size(centroids, 1);

% You need to return the following variables correctly.
idx = zeros(size(X,1), 1);
% ===================== YOUR CODE HERE =====================
%
% Note: You can use a for-loop over the examples to compute this.
%

m = size(X,1);

for i = 1:m
    distance_array = zeros(1,K);
    for j = 1:K
        distance_array(1,j) = sqrt(sum(power((X(i,:)-centroids(j,:)),2)));
    end
    [~, d_idx] = min(distance_array);
    idx(i,1) = d_idx;
end

% =========================================================

end
```

ex7_pca.m   featureNormalize.m   findClosestCentroids.m*   kMeansInitCentroids.m   pca.m   plotDataPoints.m   plotProgresskMeans.m   projectData.m

*Figure 3: Implementation of finding the closest centroids*

Command Window

```
Finding closest centroids.

Closest centroids for the first 3 examples:
 1 3 2
(the closest centroids should be 1, 3, 2 respectively)
Program paused. Press enter to continue.
```

*Figure 4: Output of the cluster center of each sample*

## Part 1.1.2: Computing centroid means

Given assignments of every point to a centroid, the second phase of the algorithm recomputes, for each centroid, the mean of the points that were assigned to it. Specially, for every centroid $k$ we set

$$\mu_k := \frac{1}{|C_k|} \sum_{i \in C_k} x^{(i)}$$

where $C_k$ is the set of examples that are assigned to centroid $k$. Concretely, if two examples say $x^{(3)}$ and $x^{(5)}$ are assigned to centroid $k = 2$, then we should update $\mu_2 = \frac{1}{2}(x^{(3)} + x^{(5)})$. We should now complete the code in **computeCentroids.m**. We can implement this function using a loop over the centroids. We can also use a loop over the examples; but if we can use a vectorized implementation that does not use such a loop, our code may run faster.

ex7.m



*Figure 5: Computing centroid means (ex7 p2)*

Once we have completed the code in **computeCentroids.m**, the code below will run our code and output the centroids after the first step of *K*-means. The centroids should be:

- [ 2.428301 3.157924 ]
- [ 5.813503 2.633656 ]
- [ 7.119387 3.616684 ]

computeCentroids.m



*Figure 6: Implementation of computing centroid means*



*Figure 7: Output of computing centroid means*

## Part 1.2: K-means on example dataset

After we have completed the two functions (findClosestCentroids and computeCentroids), the code below will run the *K*-means algorithm on a toy 2D dataset to help us understand how *K*-means works. Our functions are called from inside the **runKmeans.m** script. We encourage to take a look at the function to understand how it works. Notice that the code calls the two functions we implemented in a loop.

As before, we have to divide the points into three categories with 10 iterations. The center points of all three classes are initialized.

ex7.m



```
68
69    %% =================== Part 3: K-Means Clustering ===========================
70    %  After you have completed the two functions computeCentroids and
71    %  findClosestCentroids, you have all the necessary pieces to run the
72    %  kMeans algorithm. In this part, you will run the K-Means algorithm on
73    %  the example dataset we have provided.
74    %
75  - fprintf('\nRunning K-Means clustering on example dataset.\n\n');
76
77    % Load an example dataset
78  - load('ex7data2.mat');
79
80    % Settings for running K-Means
81  - K = 3;
82  - max_iters = 10;
83
84    % For consistency, here we set centroids to specific values
85    % but in practice you want to generate them automatically, such as by
86    % settings them to be random examples (as can be seen in
87    % kMeansInitCentroids).
88  - initial_centroids = [3 3; 6 2; 8 5];
89
90    % Run K-Means algorithm. The 'true' at the end tells our function to plot
91    % the progress of K-Means
92  - [centroids, idx] = runkMeans(X, initial_centroids, max_iters, true);
93  - fprintf('\nK-Means Done.\n\n');
94
95  - fprintf('Program paused. Press enter to continue.\n');
96  - pause;
```

*Figure 8: K-means clustering (ex7 p3)*

We have to divide the points into three categories, with 10 iterations.



*Figure 9: Before iteration*

After 10 iterations the image can see that the three piles of points are well divided into three categories. The picture also shows the movement track of the center point.



*Figure 10: Image after 10 iterations*

## Part 1.3: Random initialization

This part mainly realizes the selection of initial center points. We need to randomly select K points from the training set, and then use these points as the positions of K center points. A randomly selected code has been given in the test.

ex7.m



```
98     %% ============= Part 4: K-Means Clustering on Pixels ===============
99     %  In this exercise, you will use K-Means to compress an image. To do this,
100    %   you will first run K-Means on the colors of the pixels in the image and
101    %   then you will map each pixel on to it's closest centroid.
102    %
103    %  You should now complete the code in kMeansInitCentroids.m
104    %
105
106 -  fprintf('\nRunning K-Means clustering on pixels from an image.\n\n');
107
108    %  Load an image of a bird
109 -  A = double(imread('bird_small.png'));
110
111    % If imread does not work for you, you can try instead
112    %   load ('bird_small.mat');
113
114 -  A = A / 255; % Divide by 255 so that all values are in the range 0 - 1
115
116    % Size of the image
117 -  img_size = size(A);
118
119    % Reshape the image into an Nx3 matrix where N = number of pixels.
120    % Each row will contain the Red, Green and Blue pixel values
121    % This gives us our dataset matrix X that we will use K-Means on.
122 -  X = reshape(A, img_size(1) * img_size(2), 3);
123
124    % Run your K-Means algorithm on this data
125    % You should try different values of K and max_iters here
126 -  K = 16;
127 -  max_iters = 10;
128
```

*Figure 11: K-means clustering on pixel (ex7 p4 s-1)*

```
109 -     A = double(imread('bird_small.png'));
110
111       % If imread does not work for you, you can try instead
112       %   load ('bird_small.mat');
113
114 -     A = A / 255; % Divide by 255 so that all values are in the range 0 - 1
115
116       % Size of the image
117 -     img_size = size(A);
118
119       % Reshape the image into an Nx3 matrix where N = number of pixels.
120       % Each row will contain the Red, Green and Blue pixel values
121       % This gives us our dataset matrix X that we will use K-Means on.
122 -     X = reshape(A, img_size(1) * img_size(2), 3);
123
124       % Run your K-Means algorithm on this data
125       % You should try different values of K and max_iters here
126 -     K = 16;
127 -     max_iters = 10;
128
129       % When using K-Means, it is important the initialize the centroids
130       % randomly.
131       % You should complete the code in kMeansInitCentroids.m before proceeding
132 -     initial_centroids = kMeansInitCentroids(X, K);
133
134       % Run K-Means
135 -     [centroids, idx] = runkMeans(X, initial_centroids, max_iters);
136
137 -     fprintf('Program paused. Press enter to continue.\n');
138 -     pause;
139
```

drawLine.m | ex7.m | ex7.mlx | ex7_companion.mlx | ex7_pca.m | featureNormalize.m | findClosestCentroids.m* | kMeansInitCentroids.

*Figure 12: K-means clustering on pixel (ex7 p4 s-2)*

```matlab
1    function centroids = kMeansInitCentroids(X, K)
2    %KMEANSINITCENTROIDS This function initializes K centroids that are to be
3    %used in K-Means on the dataset X
4    %   centroids = KMEANSINITCENTROIDS(X, K) returns K initial centroids to be
5    %   used with the K-Means on the dataset X
6    %
7
8    % You should return this values correctly
9    centroids = zeros(K, size(X, 2));
10
11   % ====================== YOUR CODE HERE ======================
12   % Instructions: You should set centroids to randomly chosen examples from
13   %               the dataset X
14   %
15
16   % Initialize the centroids to be random examples
17   % Randomly reorder the indices of examples
18   randidx = randperm(size(X, 1));
19   % Take the first K examples as centroids
20   centroids = X(randidx(1:K), :);
21
22   % ============================================================
23
24   end
25
```

findClosestCentroids.m*  ×  kMeansInitCentroids.m  ×  pca.m  ×  plotDataPoints.m  ×  plotProgresskMeans.m  ×  projectData.m  ×  recoverData.m  ×

kMeansInitCentroids     Ln 1     Col 1     OVR

*Figure 13: Implementation of K-means clustering on pixel*

## Part 1.4: Image compression with K-means

In this section, K-means clustering is used for image compression. In the direct 24-bit color representation of the image, each pixel is represented as three 8-bit unsigned integers (ranging from 0 to 255), which are the intensity values of red, green, and blue. This kind of coding is also called RGB coding.

**Purpose:** Compress thousands of colors of the original image into 16 colors.

The RGB values of all pixels in the original image are used as data samples, and K-means clustering is used to cluster the points in the three-dimensional RGB space into 16 types.

ex7.m

```matlab
141     %% ================= Part 5: Image Compression =====================
142     %  In this part of the exercise, you will use the clusters of K-Means to
143 -    fprintf('\nApplying K-Means to compress an image.\n\n');
144     % Find closest cluster members
145 -    idx = findClosestCentroids(X, centroids);
146
147     % Essentially, now we have represented the image X as in terms of the
148     % indices in idx.
149
150     % We can now recover the image from the indices (idx) by mapping each pixel
151     % (specified by it's index in idx) to the centroid value
152 -    X_recovered = centroids(idx,:);
153
154     % Reshape the recovered image into proper dimensions
155 -    X_recovered = reshape(X_recovered, img_size(1), img_size(2), 3);
156
157     % Display the original image
158 -    subplot(1, 2, 1);
159 -    imagesc(A);
160 -    title('Original');
161
162     % Display compressed image side by side
163 -    subplot(1, 2, 2);
164 -    imagesc(X_recovered)
165 -    title(sprintf('Compressed, with %d colors.', K));
166
167
168 -    fprintf('Program paused. Press enter to continue.\n');
169 -    pause;
170
```

computeCentroids.m*  ×  displayData.m  ×  drawLine.m  ×  **ex7.m**  ×  ex7.mlx  ×  ex7_companion.mlx  ×  ex7_pca.m  ×  featureNormalize.m  ×  findC

*Figure 14: Image Compression (ex7 p5)*

## Part 1.4.1: K-means on pixels

After completing the previous tasks, we can directly achieve this goal. As shown in the figure below, the left is the original image, and the right is the compressed image (16 colors).

*Figure 15: Output of K-means clustering on pixels*

Use K-means for image compression. Taking a picture of $128 \times 128$ as an example, using RGB, a total of $128 \times 128 \times 24 = 393216$ bits are required. Here we compress it, divide all colors into 16 categories, and replace the colors in the entire category with the colors corresponding to the centroid. The space can be compressed to $16 \times 24 + 128 \times 128 \times 4 = 65920$ bits. Using the example provided in the title, the effect is probably as follows:



*Figure 16: Original and reconstructed image (when using K-means to compress the image)*

# Part 2: Principal Component Analysis

In this exercise, we will use principal component analysis (PCA) to perform dimensionality reduction. We will first experiment with an example 2D dataset to get intuition on how PCA works, and then use it on a bigger dataset of 5000 face image dataset. The code provided, will help us step through the first half of the exercise.

## Part 2.1: Example Dataset

First reduce the two-dimensional vector in the example to one dimension.

ex7_pca.m

```
This file uses Cell Mode. For information, see the rapid code iteration video, the publishing video, or help.
23
24      %% ================= Part 1: Load Example Dataset  ==================
25      %  We start this exercise by using a small dataset that is easily to
26      %  visualize
27      %
28 -     fprintf('Visualizing example dataset for PCA.\n\n');
29
30      %  The following command loads the dataset. You should now have the
31      %  variable X in your environment
32 -     load ('ex7data1.mat');
33
34      %  Visualize the example dataset
35 -     plot(X(:, 1), X(:, 2), 'bo');
36 -     axis([0.5 6.5 2 8]); axis square;
37
38 -     fprintf('Program paused. Press enter to continue.\n');
39 -     pause;
40
41
```

ex7.m × | ex7.mlx × | ex7_companion.mlx × | ex7_pca.m × | featureNormalize.m × | findClosestCentroids.asv × | findClosestCentr

script                                        Ln 1    Col 1    OVR

*Figure 17: Loading the example dataset (ex7_pca  p1)*

*Figure 18: Sample dataset 1*

## Part 2.2: Implementing PCA

In this part of the exercise, we will implement PCA. PCA consists of two computational steps: First, we compute the covariance matrix of the data. Then, we use MATLAB's svd function to compute the eigenvectors $U_1, U_2, \ldots, U_n$. These will correspond to the principal components of variation in the data. Before using PCA, it is important to first normalize the data by subtracting the mean value of each feature from the dataset, and scaling each dimension so that they are in the same range. In the code below, this normalization has been performed for us using the **featureNormalize function**.

After normalizing the data, we can run PCA to compute the principal components. Our task is to complete the code in **pca.m** to compute the principal components of the dataset. First, we should compute the covariance matrix of the data, which is given by:

$$\Sigma = \frac{1}{m} X^T X$$

where $X$ is the data matrix with examples in rows, and $m$ is the number of examples. Note that $\Sigma$ is a $n \times n$ matrix and not the summation operator.

ex7_pca.m

```matlab
41
42      %% =============== Part 2: Principal Component Analysis ===============
43      %   You should now implement PCA, a dimension reduction technique. You
44      %   should complete the code in pca.m
45      %
46  -    fprintf('\nRunning PCA on example dataset.\n\n');
47
48      %   Before running PCA, it is important to first normalize X
49  -    [X_norm, mu, sigma] = featureNormalize(X);
50
51      %   Run PCA
52  -    [U, S] = pca(X_norm);
53
54      %   Compute mu, the mean of the each feature
55
56      %   Draw the eigenvectors centered at mean of data. These lines show the
57      %   directions of maximum variations in the dataset.
58  -    hold on;
59  -    drawLine(mu, mu + 1.5 * S(1,1) * U(:,1)', '-k', 'LineWidth', 2);
60  -    drawLine(mu, mu + 1.5 * S(2,2) * U(:,2)', '-k', 'LineWidth', 2);
61  -    hold off;
62
63  -    fprintf('Top eigenvector: \n');
64  -    fprintf(' U(:,1) = %f %f \n', U(1,1), U(2,1));
65  -    fprintf('\n(you should expect to see -0.707107 -0.707107)\n');
66
67  -    fprintf('Program paused. Press enter to continue.\n');
68  -    pause;
```

ex7.m  ✕  ex7.mlx  ✕  ex7_companion.mlx  ✕  ex7_pca.m  ✕  featureNormalize.m  ✕  findClosestCentroids.asv  ✕  findClosestCentroids.m  ✕  kMear ▶

script                    Ln  1        Col  1      OVR

*Figure 19: Loading the example dataset (ex7_pca  p2)*

pca.m

```matlab
 1     function [U, S] = pca(X)
 2     %PCA Run principal component analysis on the dataset X
 3     %   [U, S, X] = pca(X) computes eigenvectors of the covariance matrix of X
 4     %   Returns the eigenvectors U, the eigenvalues (on diagonal) in S
 5     %
 6
 7     % Useful values
 8     [m, n] = size(X);
 9
10     % You need to return the following variables correctly.
11     U = zeros(n);
12     S = zeros(n);
13
14     % ====================== YOUR CODE HERE ======================
15     % Instructions: You should first compute the covariance matrix. Then, you
16     %               should use the "svd" function to compute the eigenvectors
17     %               and eigenvalues of the covariance matrix.
18     %
19     % Note: When computing the covariance matrix, remember to divide by m (the
20     %       number of examples).
21     %
22
23     Sigma = 1.0/m .* X' * X;
24
25     [U, S, V] = svd(Sigma);
26
27     % ==========================================================
28
29     end
```
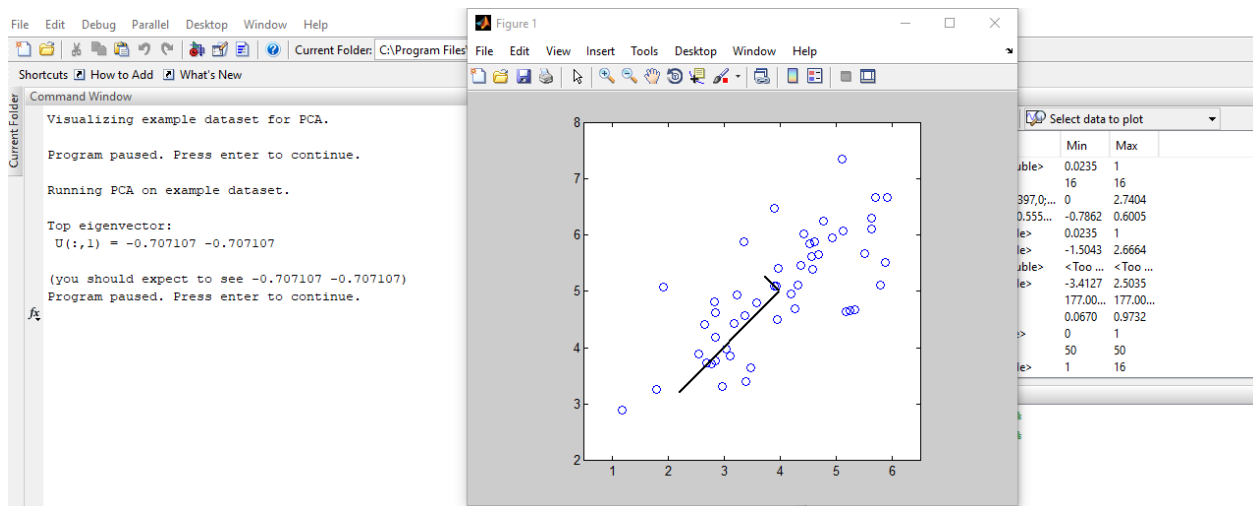
*Figure 20: Implementing PCA*



*Figure 21: The eigenvectors of the computed dataset*

## Part 2.3: Dimensionality Reduction with PCA

After obtaining the principal components, we use them to reduce the dimensionality of the data and project the data to a lower-dimensional space.

In practical applications, if we are using linear regression or neural network algorithms, we can use dimensionality reduction data to replace the original data. This can make the model train faster because the input value has a lower dimensionality.

## Part 2.3.1: Projecting the data onto the principal components

We should now complete the code in **projectData.m**. Specially, we are given a dataset X, the principal components U, and the desired number of dimensions to reduce to K. We should project each example in X onto the top K components in U. Note that the top K components in U are given by the first K columns of U, that is U_reduce = U(:, 1:K).

Once we have completed the code in **projectData.m**, run the code below to project the first example onto the first dimension and we should see a value of about 1.481 (or possibly -1.481, if you got $-U_1$ instead of $U_1$).

ex7_pca.m

```
71      %% ================== Part 3: Dimension Reduction ==================
72      %  You should now implement the projection step to map the data onto the
73      %  first k eigenvectors. The code will then plot the data in this reduced
74      %  dimensional space.  This will show you what the data looks like when
75      %  using only the corresponding eigenvectors to reconstruct it.
76      %
77      %  You should complete the code in projectData.m
78      %
79  -   fprintf('\nDimension reduction on example dataset.\n\n');
80
81      %  Plot the normalized dataset (returned from pca)
82  -   plot(X_norm(:, 1), X_norm(:, 2), 'bo');
83  -   axis([-4 3 -4 3]); axis square
84
85      %  Project the data onto K = 1 dimension
86  -   K = 1;
87  -   Z = projectData(X_norm, U, K);
88  -   fprintf('Projection of the first example: %f\n', Z(1));
89  -   fprintf('\n(this value should be about 1.481274)\n\n');
90
91  -   X_rec  = recoverData(Z, U, K);
92  -   fprintf('Approximation of the first example: %f %f\n', X_rec(1, 1), X_rec(1, 2));
93  -   fprintf('\n(this value should be about  -1.047419 -1.047419)\n\n');
94
95      %  Draw lines connecting the projected points to the original points
96  -   hold on;
97  -   plot(X_rec(:, 1), X_rec(:, 2), 'ro');
98  -   for i = 1:size(X_norm, 1)
99  -       drawLine(X_norm(i,:), X_rec(i,:), '--k', 'LineWidth', 1);
100 -   end
```

*Figure 22: Dimension reduction s-1 (ex7_pca p3)*

```
79 -     fprintf('\nDimension reduction on example dataset.\n\n');
80
81       %  Plot the normalized dataset (returned from pca)
82 -     plot(X_norm(:, 1), X_norm(:, 2), 'bo');
83 -     axis([-4 3 -4 3]); axis square
84
85       %  Project the data onto K = 1 dimension
86 -     K = 1;
87 -     Z = projectData(X_norm, U, K);
88 -     fprintf('Projection of the first example: %f\n', Z(1));
89 -     fprintf('\n(this value should be about 1.481274)\n\n');
90
91 -     X_rec  = recoverData(Z, U, K);
92 -     fprintf('Approximation of the first example: %f %f\n', X_rec(1, 1), X_rec(1, 2));
93 -     fprintf('\n(this value should be about  -1.047419 -1.047419)\n\n');
94
95       %  Draw lines connecting the projected points to the original points
96 -     hold on;
97 -     plot(X_rec(:, 1), X_rec(:, 2), 'ro');
98 -    for i = 1:size(X_norm, 1)
99 -         drawLine(X_norm(i,:), X_rec(i,:), '--k', 'LineWidth', 1);
100 -    end
101 -     hold off
102
103 -     fprintf('Program paused. Press enter to continue.\n');
104 -     pause;
105
```

*Figure 23: Dimension reduction s-2 (ex7_pca p3)*

projectData.m

```matlab
1   function Z = projectData(X, U, K)
2   %PROJECTDATA Computes the reduced data representation when projecting only
3   %on to the top k eigenvectors
4   %   Z = projectData(X, U, K) computes the projection of
5   %   the normalized inputs X into the reduced dimensional space spanned by
6   %   the first K columns of U. It returns the projected examples in Z.
7   %
8
9   % You need to return the following variables correctly.
10  Z = zeros(size(X, 1), K);
11
12  % ====================== YOUR CODE HERE ======================
13  % Instructions: Compute the projection of the data using only the top K
14  %               eigenvectors in U (first K columns).
15  %               For the i-th example X(i,:), the projection on to the k-th
16  %               eigenvector is given as follows:
17  %                    x = X(i, :)';
18  %                    projection_k = x' * U(:, k);
19  %
20  for i=1:size(X, 1),
21      for j=1:K,
22          x = X(i, :)';
23          projection_k = x' * U(:, j);
24          Z(i, j) = projection_k;
25      end
26  end
27
28  % =============================================================
29
30  end
31
```

kMeansInitCentroids.m ×  pca.m ×  plotDataPoints.m ×  plotProgresskMeans.m ×  projectData.m ×  recoverData.m ×  runkMeans.m ×  submit.m ×

projectData          Ln 31    Col 1    OVR

*Figure 24: Implementation of project data sample*

Dimension reduction on example dataset.

Projection of the first example: 1.481274

(this value should be about 1.481274)

*Figure 25: Output of project data sample*

## Part 2.3.2: Reconstructing an approximation of the data

Recovering the high dimensionality from the projected low dimensionality.

recoverData.m

```matlab
1   function X_rec = recoverData(Z, U, K)
2   %RECOVERDATA Recovers an approximation of the original data when using the
3   %projected data
4   %   X_rec = RECOVERDATA(Z, U, K) recovers an approximation the
5   %   original data that has been reduced to K dimensions. It returns the
6   %   approximate reconstruction in X_rec.
7   %
8
9   % You need to return the following variables correctly.
10  X_rec = zeros(size(Z, 1), size(U, 1));
11
12  % ====================== YOUR CODE HERE ======================
13  % Instructions: Compute the approximation of the data by projecting back
14  %               onto the original space using the top K eigenvectors in U.
15  %
16  %               For the i-th example Z(i,:), the (approximate)
17  %               recovered data for dimension j is given as follows:
18  %                    v = Z(i, :)';
19  %                    recovered_j = v' * U(j, 1:K)';
20  %               Notice that U(j, 1:K) is a row vector.
21  for i=1:size(Z, 1),
22      for j=1:size(U,1),
23          v = Z(i, :)';
24          recovered_j = v' * U(j, 1:K)';
25          X_rec(i, j) = recovered_j;
26      end
27  end
28
29  % ============================================================
30
31  end
32
```

kMeansInitCentroids.m × | pca.m × | plotDataPoints.m × | plotProgresskMeans.m × | projectData.m × | recoverData.m × | runkMeans.m × | submit.m ×

recoverData    Ln 31    Col 4    OVR

*Figure 26: Implementation recover the data by projecting*

```
Approximation of the first example: -1.047419 -1.047419

(this value should be about  -1.047419 -1.047419)

Program paused. Press enter to continue.
```

*Figure 27: Output of the data by projecting*

## Part 2.3.3: Visualizing the projections

The figure below shows how the original data is projected onto the low-dimensional hyperplane. The blue circle represents the original data, and the red circle represents the data after dimensionality reduction. The projection only effectively retains U 1 U_1 U1Information in a given direction.
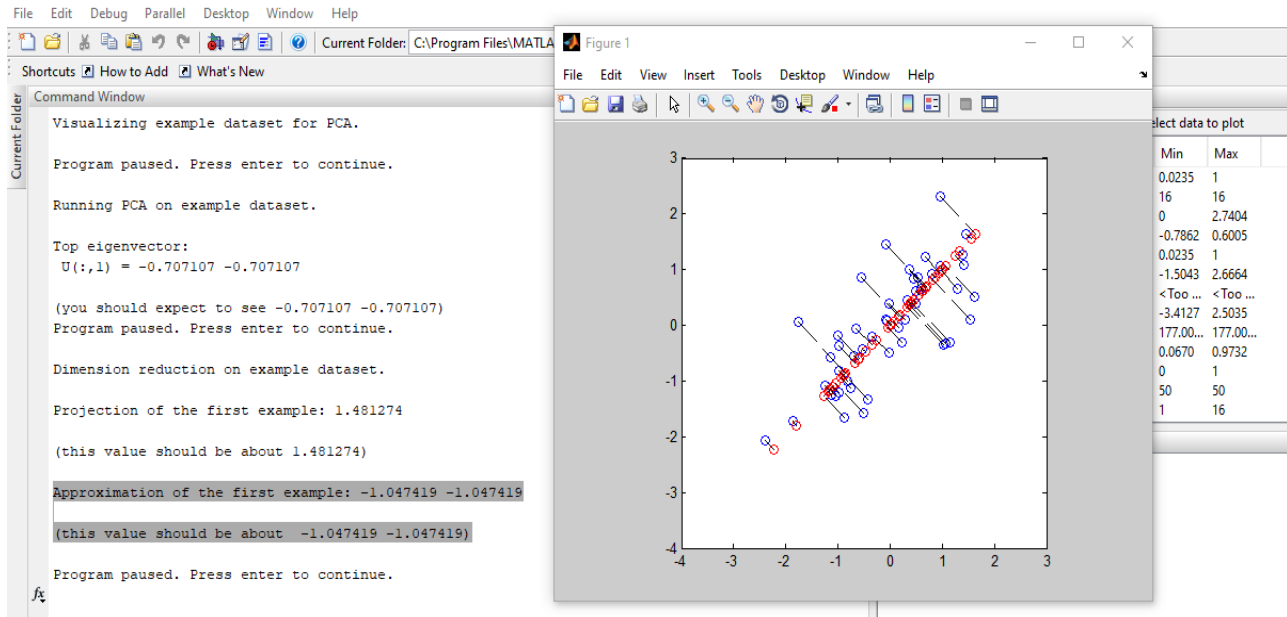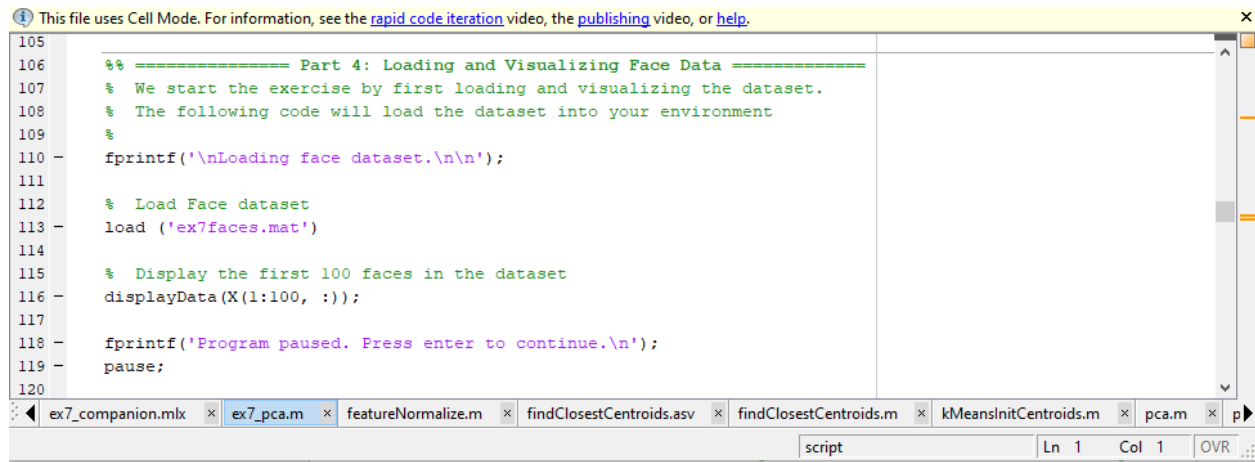


*Figure 28: PCA after normalization and projection data*

## Part 2.4: Face Image Dataset

This section we will apply PCA to face images to reduce dimensionality. Perform dimension reduction on face pictures. There are a large number of grayscale images of human faces (32/times 32) in **ex7faces.mat**, so the dimension of each vector is 32/times 32 = 1024.

ex7_pca.m

```
This file uses Cell Mode. For information, see the rapid code iteration video, the publishing video, or help.
105
106      %% =============== Part 4: Loading and Visualizing Face Data ==============
107      %  We start the exercise by first loading and visualizing the dataset.
108      %  The following code will load the dataset into your environment
109      %
110 -     fprintf('\nLoading face dataset.\n\n');
111
112      %  Load Face dataset
113 -     load ('ex7faces.mat')
114
115      %  Display the first 100 faces in the dataset
116 -     displayData(X(1:100, :));
117
118 -     fprintf('Program paused. Press enter to continue.\n');
119 -     pause;
120
ex7_companion.mlx  ×  ex7_pca.m  ×  featureNormalize.m  ×  findClosestCentroids.asv  ×  findClosestCentroids.m  ×  kMeansInitCentroids.m  ×  pca.m  ×  p
                                                          script          Ln 1    Col 1    OVR
```
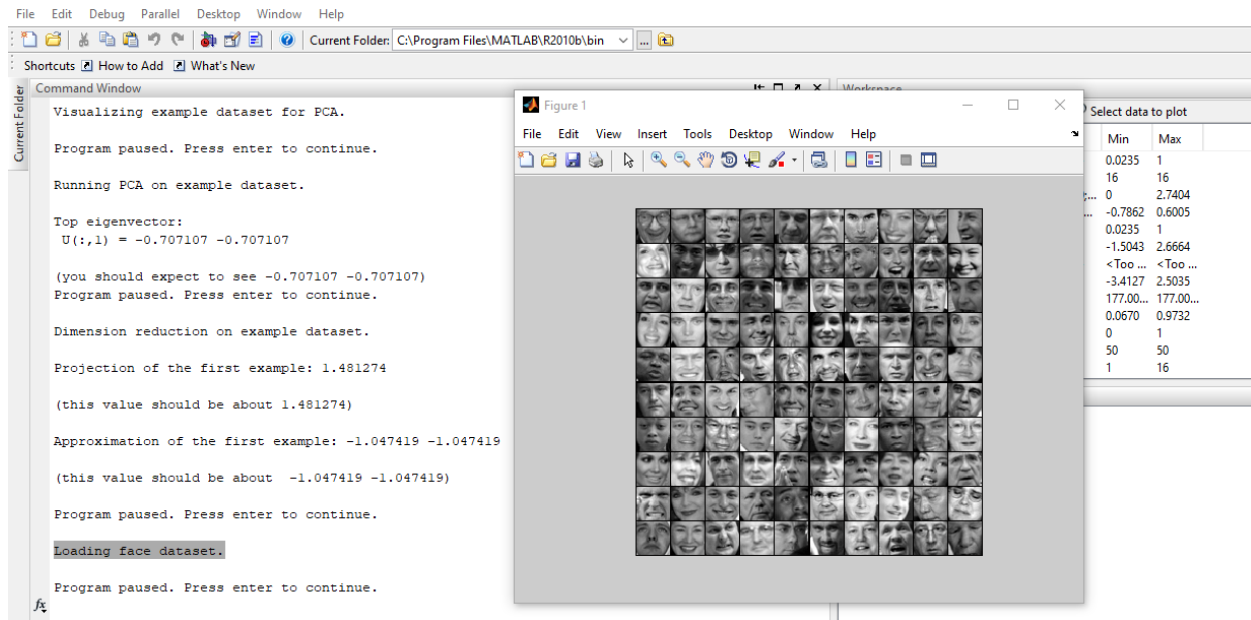
*Figure 29: Visualizing face dataset (ex7_pca  p4)*

*Figure 30: Face data sets*

## Part 2.4.1: PCA on Faces

We can do this by reshaping each of them into 32×32 Matrix to visualize these principal components, the matrix corresponds to the pixels in the original dataset. Script **ex7_pca.m** Show the front that describes the biggest change 36 A principal component. If necessary, we can also change the code to show more principal components, to see how they capture more and more details.

This file uses Cell Mode. For information, see the rapid code iteration video, the publishing video, or help.                                                                 ×

```matlab
120
121     %% ============ Part 5: PCA on Face Data: Eigenfaces  ====================
122     %  Run PCA and visualize the eigenvectors which are in this case eigenfaces
123     %  We display the first 36 eigenfaces.
124     %
125 -   fprintf(['\nRunning PCA on face dataset.\n' ...
126               '(this mght take a minute or two ...)\n\n']);
127
128     %  Before running PCA, it is important to first normalize X by subtracting
129     %  the mean value from each feature
130 -   [X_norm, mu, sigma] = featureNormalize(X);
131
132     %  Run PCA
133 -   [U, S] = pca(X_norm);
134
135     %  Visualize the top 36 eigenvectors found
136 -   displayData(U(:, 1:36)');
137
138 -   fprintf('Program paused. Press enter to continue.\n');
139 -   pause;
140
```

| ex7_companion.mlx × | ex7_pca.m × | featureNormalize.m × | findClosestCentroids.asv × | findClosestCentroids.m × | kMeansInitCentroids.m × | pca.m × | p |

script                                             Ln 1      Col 1     OVR

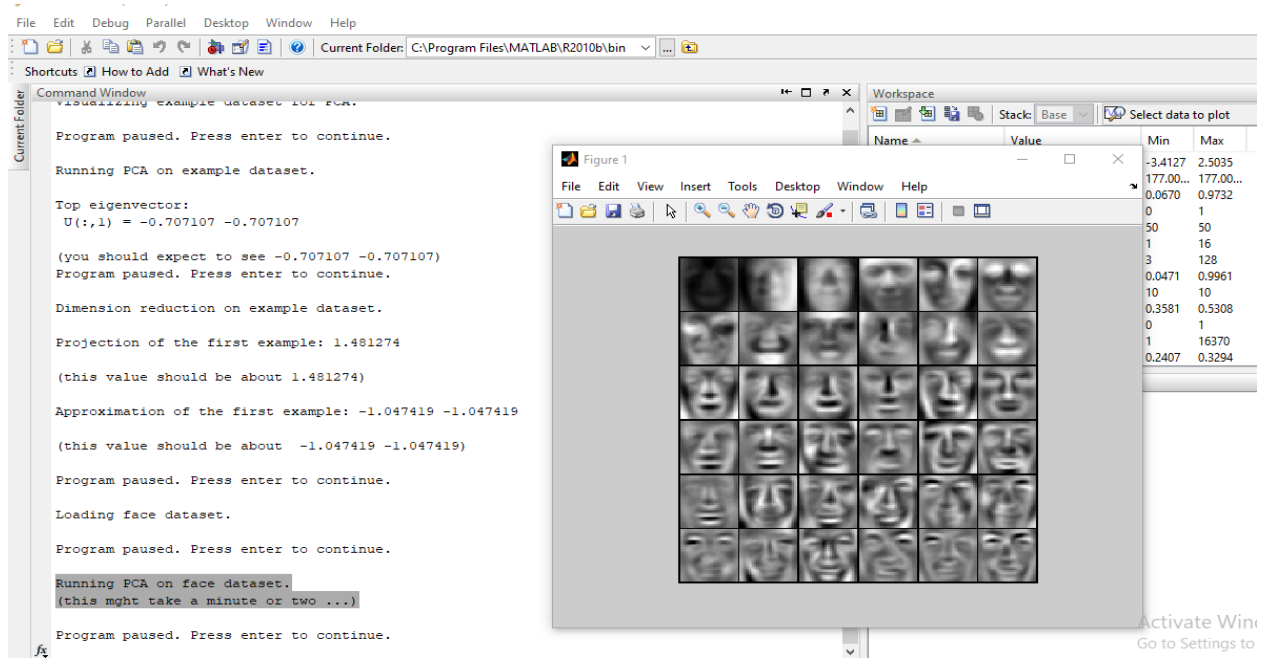*Figure 31: PCA on face dataset (ex7_pca  p5)*

*Figure 32: The principal components on the face dataset*

## Part 2.4.2: Dimensionality Reduction

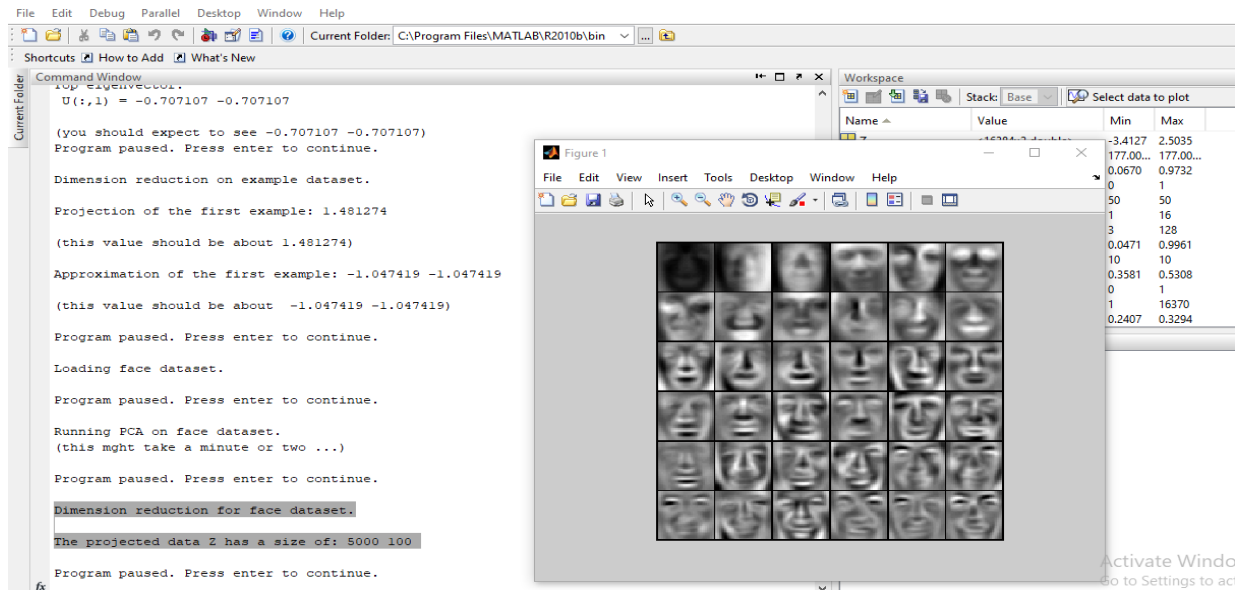Take the first 100 feature vectors for projection



*Figure 33: Principal components on the face dataset*

It can be seen that after the dimensionality is reduced, the general frame of the face is still retained, but some details are lost. The inspiration for us is that when we are training face recognition with neural networks, we can sometimes use this method to increase the speed

ex7_pca.m

```
156    %% ==== Part 7: Visualization of Faces after PCA Dimension Reduction ====
157    %  Project images to the eigen space using the top K eigen vectors and
158    %  visualize only using those K dimensions
159    %  Compare to the original input, which is also displayed
160
161    fprintf('\nVisualizing the projected (reduced dimension) faces.\n\n');
162
163    K = 100;
164    X_rec  = recoverData(Z, U, K);
165
166    % Display normalized data
167    subplot(1, 2, 1);
168    displayData(X_norm(1:100,:));
169    title('Original faces');
170    axis square;
171
172    % Display reconstructed data from only k eigenfaces
173    subplot(1, 2, 2);
174    displayData(X_rec(1:100,:));
175    title('Recovered faces');
176    axis square;
177
178    fprintf('Program paused. Press enter to continue.\n');
179    pause;
180
```

ex7_companion.mlx × | ex7_pca.m × | featureNormalize.m × | findClosestCentroids.asv × | findClosestCentroids.m × | kMeansInitCentroids.m ▶

script          Ln  1     Col  1     OVR

*Figure 34: Visualization of Faces after PCA Dimension Reduction*

The image on the left is the original image, and the image on the right is the image reconstructed after dimensionality reduction.
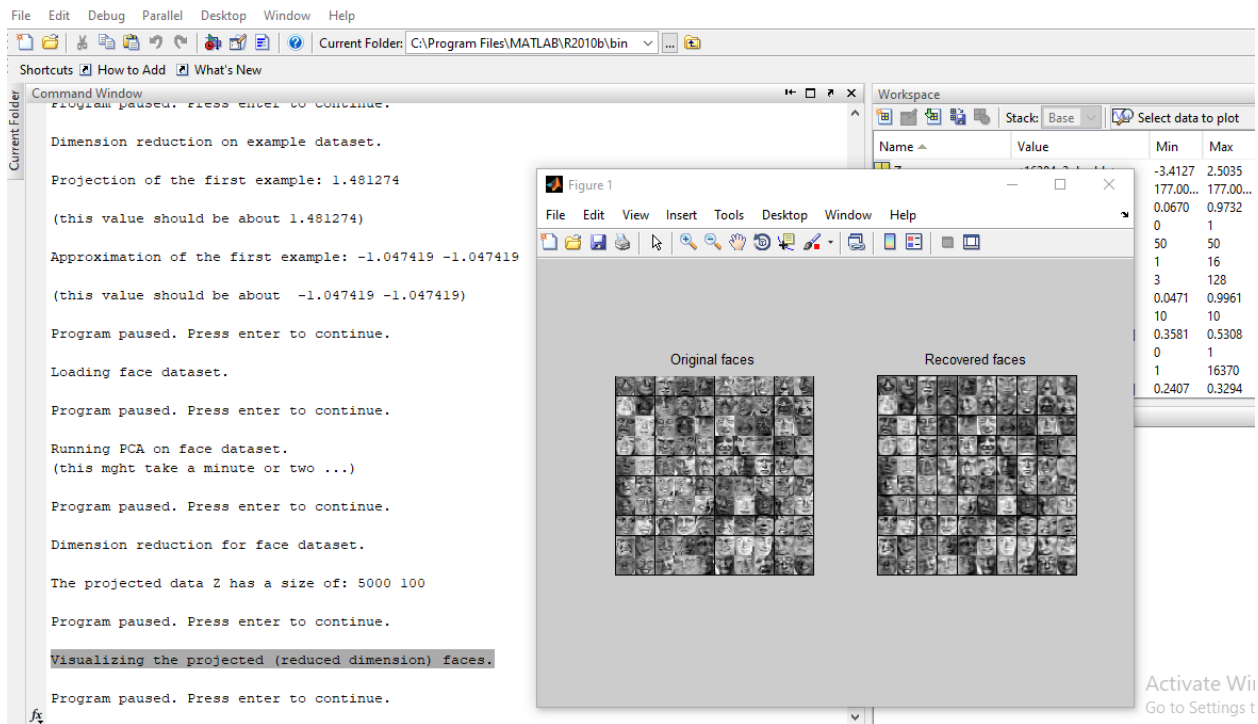
*Figure 35: Original images of faces and ones reconstructed from*

*only the top 100 principal components.*

It can be seen that after reducing the dimensions, the general framework of the human face is still preserved, but some details are lost. This inspired us that when we are training facial recognition with neural networks, we can sometimes use this method to increase speed.

## Part 2.5: PCA for visualization

PCA Commonly used for high-dimensional vector display. We will use PCA implementation to 3D data, reduce it to a 2D space, and visualize the result as a 2D scatter plot. The PCA projection can be thought of as a rotation, which selects the view that maximizes the data, which usually corresponds to the "best" view.

ex7_pca.m

```
182    %% === Part 8(a): Optional (ungraded) Exercise: PCA for Visualization ===
183    %  One useful application of PCA is to use it to visualize high-dimensional
184    %  data. In the last K-Means exercise you ran K-Means on 3-dimensional
185    %  pixel colors of an image. We first visualize this output in 3D, and then
186    %  apply PCA to obtain a visualization in 2D.
187
188 -  close all; close all; clc
189
190    % Re-load the image from the previous exercise and run K-Means on it
191    % For this to work, you need to complete the K-Means assignment first
192 -  A = double(imread('bird_small.png'));
193
194    % If imread does not work for you, you can try instead
195    %   load ('bird_small.mat');
196
197 -  A = A / 255;
198 -  img_size = size(A);
199 -  X = reshape(A, img_size(1) * img_size(2), 3);
200 -  K = 16;
201 -  max_iters = 10;
202 -  initial_centroids = kMeansInitCentroids(X, K);
203 -  [centroids, idx] = runkMeans(X, initial_centroids, max_iters);
204
205    %  Sample 1000 random indexes (since working with all the data is
206    %  too expensive. If you have a fast computer, you may increase this.
207 -  sel = floor(rand(1000, 1) * size(X, 1)) + 1;
208
209    %  Setup Color Palette
210 -  palette = hsv(K);
211 -  colors = palette(idx(sel), :);
```

ex7_companion.mlx  ×  | ex7_pca.m  ×  | featureNormalize.m  ×  | findClosestCentroids.asv  ×  | findClosestCentroids.m  ×  | kMeansInitCentroids.m

*Figure 36: PCA for visualizing s-1 (ex7_pca p8)*

```
207 -    sel = floor(rand(1000, 1) * size(X, 1)) + 1;
208
209      %  Setup Color Palette
210 -    palette = hsv(K);
211 -    colors = palette(idx(sel), :);
212
213      %  Visualize the data and centroid memberships in 3D
214 -    figure;
215 -    scatter3(X(sel, 1), X(sel, 2), X(sel, 3), 10, colors);
216 -    title('Pixel dataset plotted in 3D. Color shows centroid memberships');
217 -    fprintf('Program paused. Press enter to continue.\n');
218 -    pause;
219
220      %% === Part 8(b): Optional (ungraded) Exercise: PCA for Visualization ===
221      % Use PCA to project this cloud to 2D for visualization
222
223      % Subtract the mean to use PCA
224 -    [X_norm, mu, sigma] = featureNormalize(X);
225
226      % PCA and project the data to 2D
227 -    [U, S] = pca(X_norm);
228 -    Z = projectData(X_norm, U, 2);
229
230      % Plot in 2D
231 -    figure;
232 -    plotDataPoints(Z(sel, :), idx(sel), K);
233 -    title('Pixel dataset plotted in 2D, using PCA for dimensionality reduction');
234 -    fprintf('Program paused. Press enter to continue.\n');
235 -    pause;
236
```

ex7_companion.mlx   ×   ex7_pca.m   ×   featureNormalize.m   ×   findClosestCentroids.asv   ×   findClosestCentroids.m   ×   kMeansInitCentroids.m   ×

script                              Ln  1      Col  1      OVR

*Figure 37: PCA for visualizing s-2 (ex7_pca p8)*

**PCA** Commonly used for high-dimensional vector display. As shown below, use **K-means** to classify points in three-dimensional space.
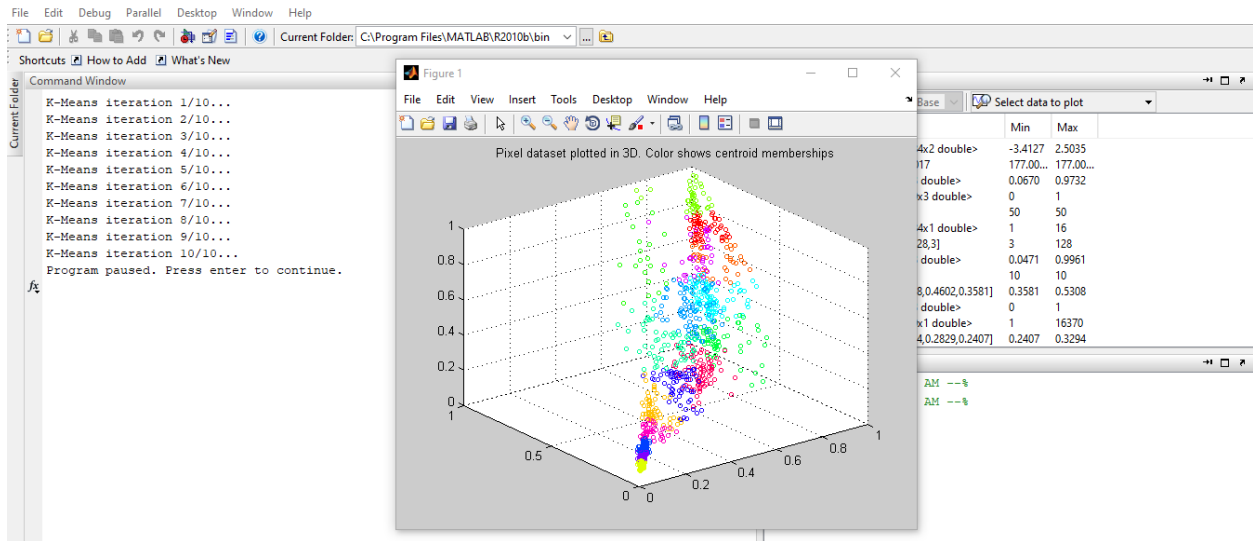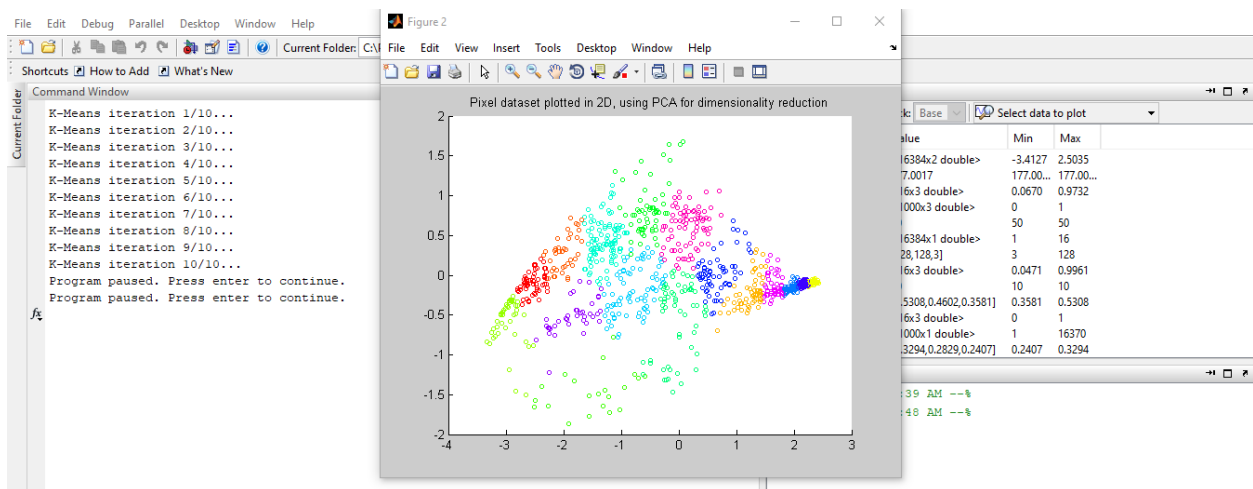
*Figure 38: Raw data in dimensions*



*Figure 39: Use PCA Produced 2D visualization*