### 24-787 Artificial Intelligence and Machine Learning for Engineering Design

### Dimensionality Reduction

**Naming and Folder structure:** Name a root folder "andrewid_hw#". Example: lkara_hw1. In that root folder, create subfolders such as q1, q2, q3 etc. each corresponding to a question in the assignment. The files and documents related to each question should go into the corresponding subfolder.

**File types:** Your folders and subfolders should not contain any files except .pdf, .doc, .docx, .m, .mat, .txt .zip files. If you want to take photos of your assignment, just make sure that combine all the jpg into a single pdf file and make it clear. Unorganized and illegible files will be penalized. Take care in arranging your illustrations, written solutions, photos. If any, do not scan the your hand-written solutions in the highest resolution.

**What to include:** Only submit the required files that you create or modify. For programming questions, include your source code (.m and/or .mat file rather than text) in your submission. All other files, including the ones we provide as supporting functionality, are unnecessary because we already have copies of them. Unless we tell you otherwise, make sure that you only submit the file you edited/changed. For .m files, we expect to click the "run" button and everything should work (obviously, after we include the additional support files and data that was given to you with the assignment).

**Readme files:** Provide a readme.txt within each questions' subfolders when necessary to aid grading. This is useful when there are many files of the same type. If your folder structure and your file names are logical, you should not need a readme.txt file.

**Submission:** Zip the entire assignment by compressing the root folder. Example: lkara_hw1.zip. Do not use .rar**.** Make sure your entire submission file is less than 20MB when zipped. You shall use the online homework submission platform. The submission entry can be found at the top right of homework/assignment page. Make sure your submission is in .zip format, otherwise it cannot be uploaded. Feel free to re-submit your homework before deadline if you find any mistake in your previous submission. Only the last submission will be graded.

1. *(25 points)*

   Suppose you are given the following two-dimensional dataset (see *q1data.txt*):

   $$X = \begin{bmatrix} 2.5 & 2.4 \\ 0.5 & 0.7 \\ 2.2 & 2.9 \\ 1.9 & 2.2 \\ 3.1 & 3.0 \\ 2.3 & 2.7 \\ 2.0 & 1.6 \\ 1.0 & 1.1 \\ 1.5 & 1.6 \\ 1.1 & 0.9 \end{bmatrix}$$

   (a) Compute the first and second principal components ($e_1$ and $e_2$) of this dataset. Show all your work. Plot the data in the original x1-x2 space and show the principal components.

   (b) Transform the data using both principal components (*i.e.* compute $a_1$ and $a_2$ for each data point) and plot this new representation on the a1-a2 plane.

   (c) What is the PCA-optimal one-dimensional representation of the data? In this reduced dimension, what is the range of the data (the distance between the minimum and maximum points)?

2. *(35 points)*

   Now, suppose your data is as follows (see *q2data.txt*):

   $$X = \begin{bmatrix} -2 & 1 & 2 & -3 & 4 & 1 & 0 & 3 & 0 & 2 & 1 & 1 & 2 & 3 & -2 & -3 & 2 & 1 & 0 \\ 1 & 2 & -4 & 2 & -4 & 2 & 5 & 2 & 2 & 1 & -3 & 0 & 0 & 1 & -2 & 1 & 1 & -3 & -2 \\ 1 & -3 & 2 & 1 & 0 & -3 & -5 & -1 & 3 & 3 & -2 & -3 & -2 & -1 & 1 & 0 & 5 & 4 & 2 \\ 3 & -1 & 0 & 2 & 2 & -5 & -4 & -1 & 2 & -1 & 3 & 4 & 4 & 2 & 1 & 2 & -2 & 1 & -1 \end{bmatrix}$$

   In this case, the number of samples (4) is much less than the number of dimensions (19), so you will need to use an approach similar to the face recognition problem discussed in lecture.

   (a) Compute the 4x4 inner product matrix and output all non-zero eigenvectors. How many such vectors are there? Show all steps in your derivation.

   (b) Calculate the projection of each data point onto the 3D reduced space. Show all steps in your derivation. Report the 4X3 matrix.

(c) Determine the mean-squared error between each original sample and its reconstructed version using three eigenvectors. Note that the MSE should be computed in the original 4 dimensional space.

(d) Repeat part (c), but this time assume the dimensionality of the data is reduced to two instead of three.

(e) What is the Euclidean distance between the new data vector given below and each of the four samples in the reduced three-dimensional space? Which of the four samples is most similar to this new vector?

$$Y = \begin{bmatrix} 1 & 3 & 0 & 3 & -2 & 2 & 4 & 1 & 3 & 0 & -2 & 0 & 1 & 1 & -3 & 0 & 1 & -2 & -3 \end{bmatrix}$$

(f) Perform the same analysis as part (e) (*i.e.* nearest-neighbor classifier), but this time in the original dimension space. Do the results match? Does this make intuitive sense? Why or why not?

*Hint:* Make sure your eigenvectors are always unit vectors.

3. *(40 points)*

In this programming exercise, you will be developing a simple classifier to recognize human faces. The image dataset (Yale face database) is provided with the assignment. It consists of 165 grayscale images (15 subjects, 11 images each), each of size 231x195. The provided MATLAB variables include:

| | |
|---|---|
| X | The data, encoded as a 45045x165 matrix, in which the rows represent dimensions (pixels) and the columns indicate examples. To view an image, simply reshape the appropriate column vector into a 231x195 array. For example, the following line of code displays the 142$^{nd}$ image in the dataset: `imshow(reshape(data(:,142),231,195))` |
| Y | The labels, given as a 165x1 vector. Each label takes an integer value from 1-15, indicating which subject is shown in the corresponding image. |
| testimages | A vector of indices indicating which images to withhold for testing. For example, if the number 10 is included in this vector, do not use the 10$^{th}$ image when computing principal components. |
| trainimages | The complement of `testimages`, provided for completeness. |

(a) Reduce the dimensionality of the data using PCA such that 90% of the variance is preserved. That is, find the number of eigenvectors E such that the sum of the first E eigenvalues accounts for 90% of the total sum of all eigenvalues. What is the new number of dimensions E? In your report, provide visualizations of the first five eigenfaces (those with the highest eigenvalues) and include a bar graph showing the percentage of variance explained by each principal component in descending order.

(b) Try reconstructing the first training example with varying number of eigenfaces. In your report, juxtapose the original image with its reconstructed versions using [10, 20, 30, 40, 50] eigenfaces.

(c) Implement the nearest neighbor algorithm in the reduced dimension space to classify each of the test images. Report the overall test accuracy.

(d) Repeat part (c), but in the original dimension space. Does the accuracy change? Is this the result you would expect? Why or why not?

4. **Extra credit:** *(10 points max)* Use Linear Discriminant Analysis (LDA) to classify the test images in Problem 3. You will have to study LDA on your own. Compare your accuracy to the result in Problem 3, part (c).