

## AMATH 482 HW 4

### Use of the Singular Value Decomposition in Classification

**Zachary McNulty**  
zmcnulty, ID: 1636402

**Abstract:** In this paper, we will explore the power and limitations of the SVD in extracting fundamental structure from data. To do so, we will attempt to form a low-rank approximation of the human face and explore the consequences of using data that has not been preprocessed for SVD. Furthermore, we will use these ideas in order to develop a coordinate system ideal for classification via Support Vector Machines (SVM), k-Nearest Neighbors (KNN), and Decision Tree Learning (DTL). As an example, we will run this algorithm on a series of music samples from different bands and genres and analyze their performance.

# 1 Introduction and Overview

The SVD is a powerful tool for understanding the fundamental structure of data, but it has certain limitations. One of these limitations is that it does not handle translated data well. Because the SVD relies on simple linear combinations of its principal modes to capture features in the data, shifted data measurements cannot be well-captured by their unshifted counterparts. This can lead to an artificially high approximation of the dimension of the system, reducing the low-rank benefits SVD provides. To highlight this shortcoming, we will perform SVD analysis on two groups of images from the Yale Faces Database: preprocessed, cropped headshots and uncropped, raw headshots. Later on in this report, we will use what we learned about the SVD to develop powerful classification techniques. Specifically, we will use SVD to find a low-dimension representation of our data and choose some of these dimensions, those which show the most variability between the groups we are trying to classify, to form a coordinate system for classification. Once we have developed this coordinate system, we will use Support Vector Machines (SVM), k-Nearest Neighbors (KNN), and Decision Tree Learning (DCL) to perform the classification, analyzing the results of each of these three models.

## 2 Theoretical Background

The Singular Value Decomposition is a diagonalization of a given matrix  $A$  that focuses on the rotations/reflections and stretching a vector undergoes when transformed by  $A$ . Formally, for any matrix  $A \in \mathbb{C}^{m \times n}$  there exists a diagonalization of the form:

$$A = U\Sigma V^* \quad (1)$$

where  $U \in \mathbb{C}^{m \times m}$  and  $V \in \mathbb{C}^{n \times n}$  are unitary matrices,  $\Sigma \in \mathbb{R}_{\geq 0}^{m \times n}$  is a diagonal matrix, and  $V^*$  represents the complex conjugate transpose of  $V$ . As  $U$  and  $V$  are unitary matrices their columns each form an orthonormal basis for their respective space and the matrices have the convenient property that  $UU^* = U^*U = I_m$  and  $VV^* = V^*V = I_n$ . The diagonal entries of  $\Sigma$  are called the **singular values** of  $A$ . Intuitively, this breaks down the transformation  $A$  into three fundamental transformations: a rotation/reflection  $V^*$  within the the domain of  $A$ , a scaling of the components of the space  $X$  by  $\Sigma$ , and another rotation/reflection of the space within the codomain of  $A$ . To understand the importance of the SVD, first consider the covariance matrix for a data matrix  $X$ :

$$C_X = \frac{1}{n-1}XX^T \quad (2)$$

In this matrix, entry  $(i, j)$  corresponds to the covariance between row  $i$  and row  $j$  of  $X$  (data measurement  $i$  and  $j$ ). If  $i \neq j$ , then a high value in entry  $(i, j)$  suggests the two data measurements vary in similar ways and are thus likely redundant measurements of the same feature in the data. Consequentially, a low value at  $(i, j)$  suggests the two data measurements are fairly independent. If  $i = j$ , the diagonal of  $C_X$ , then entry  $(i, j)$  describes the variance in data measurement  $i$ . Generally, data measurements with a high variance are assumed to capture important features of the data while those with low variance do not. To see the importance of the SVD, consider the covariance matrix for  $Y = U^*X = \Sigma V^*$  where  $U, \Sigma, V$  are from the SVD of  $X$ :

$$C_Y = \frac{1}{n-1}YY^T = \frac{1}{n-1}\Sigma^2 \quad (3)$$

As  $C_Y$  is diagonal, all its off diagonal entries are zero: projecting the data  $X$  onto  $U^*$  produces a transformed data set  $\Sigma V^*$  with completely independent measurements, eliminating redundancy in the data. Similarly, if we chose  $Y = XV = U\Sigma$  we find the same covariance matrix, showing us  $U\Sigma$  also forms an independent set. As the columns of  $U$  and  $V$  form bases for the codomain and domain of  $X$  respectively, each give an ideal, independent coordinate system for each space. If the data samples in  $X$  are stored as rows,  $V$  gives the ideal coordinate system and  $U$  the coordinates of each sample in that system. Else, if the samples are stored as columns then  $U$  gives the ideal coordinate system and  $V$  the coordinates of each sample in that system. Furthermore, the corresponding singular value  $\sigma$  in the diagonal of  $\Sigma$  ranks the importance of each direction in this coordinate system based on the data's variance along that direction. Note for any given matrix  $X$  we can reconstruct it given these coordinates:

$$X = \sum_{i=1}^{\min(m,n)} u_i \sigma_i v_i^* \quad (4)$$

where  $u_j, \sigma_j, v_j$  are the  $j$ th column of  $U, \Sigma, V$  respectively. Since the singular values  $\sigma_j$  are a non-increasing sequence, we see each successive term is increasingly less important to the structure of  $X$ . Thus, if the singular values decrease rapidly, we can get a good approximation of  $X$  with only the first few terms of this sum. This is a low-rank approximation of the matrix  $X$ , a simplified form of the dynamics of the system that well-captures its behavior.

As we will show later, the SVD is very sensitive to translations. In this regard, we can use tools such as the Fourier Transform to move the data into a space where this translation is not as relevant. For example, if we were to have two identical sinusoidal signals shifted completely out of phase, SVD on the time series data would give two prominent modes. However, performing the SVD in the frequency domain on the Fourier Transforms of these signals will yield a single prominent mode: the frequency of the system. The Fourier Transform is the primary tool for converting between the spatial and frequency domains of a signal. Below is the transform and its inverse:

$$F(k) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-ikx} f(x) dx \quad (5)$$

$$f(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{ikx} F(k) dk \quad (6)$$

Where  $k$  is the wavenumber. It may not be clear from **Equations 5** or **6** why this encodes frequency information, but expanding  $e^{-ikx}$  provides some insight:

$$e^{-ikx} = \cos(kx) - i \sin(kx) \quad (7)$$

**Equation 7** makes it clear why then  $k$  represents the wavenumber. Thus the transform is essentially taking the inner product of our function with a bunch of periodic sines and cosines where, roughly speaking, frequencies/wavenumbers more similar to our signal get assigned higher values. In this way, spatial information is converted to the frequency domain.

Simplifying the system allows us to extract the important characteristics of a system. Not only does this reduce the amount of computational work required to work with the system, but it lowers the number of parameters involved and makes overfitting a training algorithm less likely. The three supervised learning algorithms used in this paper are KNN, SVM, and DTL. Supervised learning follows a simple idea: given a set of labeled data measurements from a set of classes, learn to classify future data measurements into one of the pre-defined classes. KNN simply classifies new points by finding which of the labeled points the new point is closest to (or more generally, of the  $K$  nearest labeled points, which class is most common). SVM is essentially an optimization problem that tries to find the an optimal hyperplane  $\vec{w}x + b = 0$  where all points  $x$  satisfying  $\vec{w}x + b < 0$  are labeled as class one and all points satisfying  $\vec{w}x + b > 0$  are labeled as class two. In this case, the optimal hyperplane is defined as the one in which there are as few errors as possible in classification on the training data and that distance separating the classes around the boundary is as large as possible. The latter constraint tries to make it less likely the variability in the input will cause a change in classification. Lastly, DTL follows a simple binary approach. For each variable in our system, scan over all possible values of that variable to be use as a separator of that data. Find the variable and its associated value that best divides the current collection of data into their respective groups. Divide the data according to this variable/value and repeat the process for both halves of the data, continuing for a set number of iterations or until some tolerance is reached. In the end, the algorithm yields a series of conditions that can be used to classify the data by following the path down the decision tree. As the success of all three of these supervised learning algorithms depends on how well they can separate the classes within the coordinate system they are working in, SVD can be very important for helping develop the proper space to work in.

### 3 Algorithm Implementation and Development

#### 3.1 Yale Faces

1. Loop over all the folders containing photos for each subject (Appendix B: code 0-31).
2. For each folder/subject, take each photo, reshape it into a column vector, and place it as a column in our data matrix  $F$ . Additionally, create an average of all the images for each subject and store it as a column in  $F_{ave}$  (Appendix B: code 35-50).

3. Mean subtract from data matrices  $F$  and  $F_{ave}$  and perform SVD (Appendix B: code 59-71).
4. Plot the singular values, analyze the principal components (columns of  $U$ ), and choose a rank  $r$  that captures at least 90% of the energy of the system (Appendix B: code 97-139).
5. Perform the rank  $r$  approximation of the faces and analyze the effectiveness of the approximation (Appendix B: code 145-163).
6. Repeat for the uncropped images and compare results (Appendix B: code 164-307).

### 3.2 Music Classification

1. Loop through the folders for each group (Appendix B: code 0-42).
2. For each group, loop through all the songs and take 5-second samples with a random starting point, storing which group the samples are from. We ignore the first 30 seconds of a song to avoid sampling the period of time before the music begins (Appendix B: code 49-116).
3. Perform the FFT on each sample to convert to frequency domain. Place half these samples into the training data set and the other half into the validation dataset (Appendix B: code 81-92).
4. Mean subtract from each sample and perform the SVD on the entire training dataset (Appendix B: code 133-134).
5. Plot singular values and determine principal components that adequately separate the data to develop a better coordinate system for classification. Choose a rank that captures at least 90% of the energy of the system (Appendix B: code 137-158).
6. Train SVM, KNN, and DTL models using training dataset, training the model on the coordinates of each of these pieces in the new coordinate system determined by SVD in the previous steps (coordinates of each music piece with respect to the basis  $U$ ) (Appendix B: code 182-203).
7. Check accuracy of the model generated by projecting each piece of data in the validation dataset onto the basis/coordinate system determined by SVD and check the model's prediction (Appendix B: code 213-226).
8. Plot classifications and analyze where misclassifications occurred (Appendix B: code 233-242).
9. Repeat this process for all three test cases.

## 4 Computational Results

### 4.1 Yale Faces

Since we stored each image as a column in our data matrix  $F$ , the SVD has a standard interpretation. As  $U$  forms an orthonormal basis for the codomain of our data matrix  $F$  and the codomain is the space where each image exists, the  $U$  matrix is our coordinate system for "face-space". We can see this when we plot the columns of  $U$  which make up the basis of face-space. As we see in **Figure 1**, each principal component captures some of the general features of a face. Similarly, just by considering the matrix multiplication involved in  $F = U\Sigma V^*$ , we see that after scaling these principal components by  $\Sigma$ , each column of  $V^*$  gives the coordinates of each image, a column in  $F$ , in this face-space defined by  $U$ : Each face in  $F$  is a linear combination of the columns of  $U$  weighted by the columns of  $\Sigma V^*$ .

When we are working with the cropped images, we see that the averaging works quite well. As all the faces are aligned in the frame, averaging them serves to balance out the lighting in the scene. Furthermore, as we see in **Figure 2**, averaging allows us to drastically reduce the size of our system. While it appears that both systems have similar trends in the singular values, it takes 1181 modes to capture 90% of energy in the system versus only 29 modes in the averaged-face system. By reducing the dimension of our system, we can greatly improve the amount of computational work we can accomplish with it and reduce any overfitting that might come from

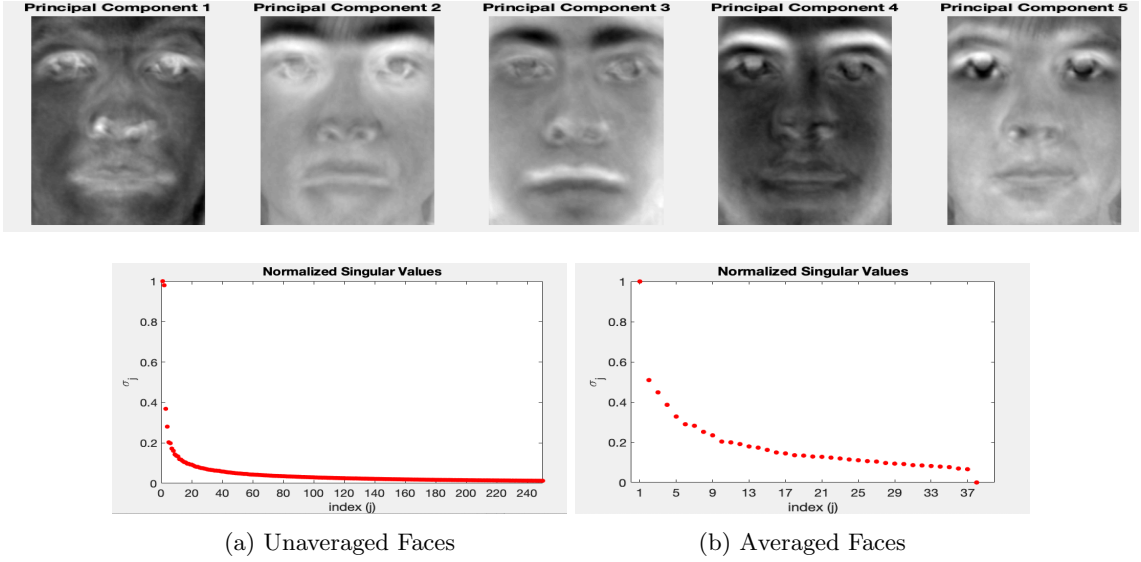


Figure 1: Face-space for the cropped, averaged images (top). The columns of  $U$  reshaped into the dimensions of the original images and plotted as an image. We can see each principal component captures some of the general features of a face. We can see there are many large singular values, suggesting face-space is fairly high-dimensional.

trying to extend this system, two key features whose importance will become clear in the Music Classification section. Conversely, with the uncropped images the averaged image does not make much sense. Since the subjects are in different locations, averaging the image introduces a large amount of blur and makes it difficult to distinguish key features of the face. Even though this reduces the dimensionality of our system, it is not clear the reduced system generated will be on any benefit. As we see in **Figure 3**, the face-space for the averaged, uncropped images is poorly recognizable. Even without averaging, the varying location of the subject makes it difficult for the SVD to pick out important features and we see a similar face-space. While the averaging in the uncropped face-space still allows us to reduce dimensionality, its cost is clearly seen in the singular values of **Figure 3**: there are many more relevant (high variance) modes than in the uncropped case. Thus, our low-rank approximations will not be as good and we lose some of the aforementioned benefits of our reduced system.

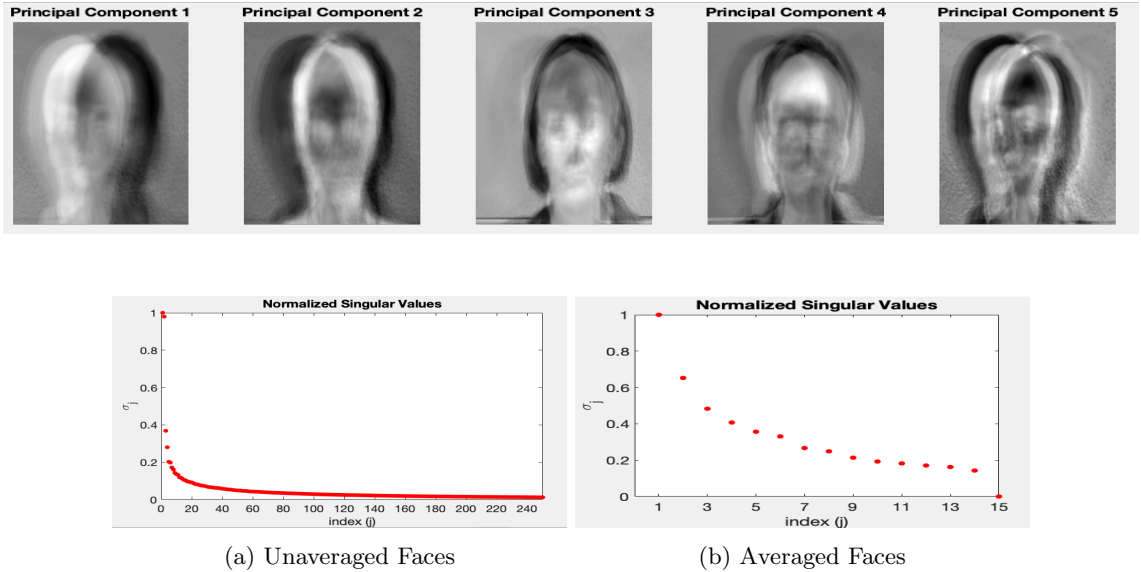


Figure 2: Face-space for the uncropped, averaged images (top). The columns of  $U$  reshaped into the dimensions of the original images and plotted as an image. We can see with the uncropped photos, the averaged image poorly captures features of face-space. This is also reflected in an increase in the decay rate of our singular values (bottom), clearly apparent in the averaged case.

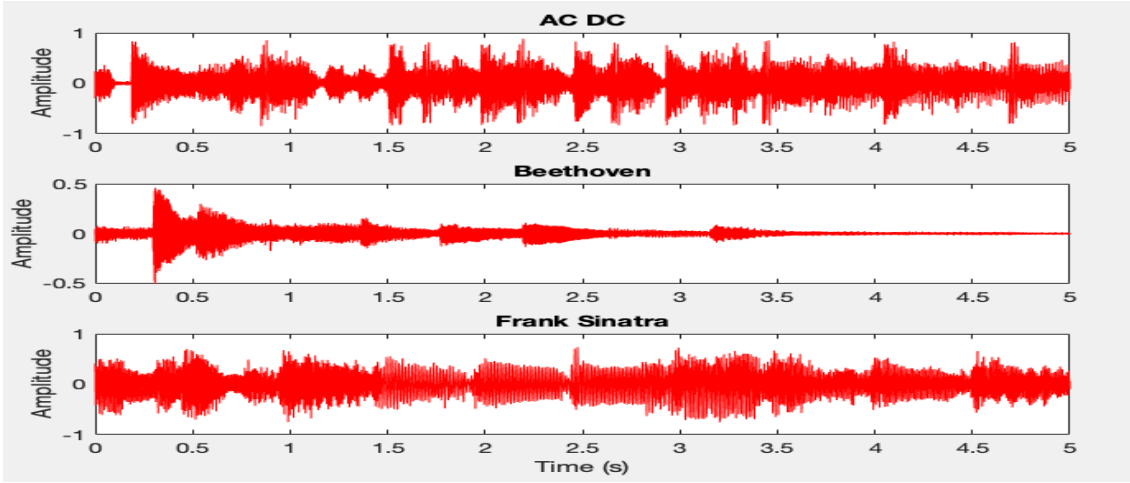


Figure 3: The music from different songs, bands, and genres each have a unique frequency pattern. As we see above, AC/DC has rapid frequency patterns which change suddenly, while Frank Sinatra and Beethoven have much more even patterns.

## 4.2 Music Classification

As we can see in **Figure 4**, each song, band, and genre have characteristic patterns that make them what they are. Using SVD, we can break down these patterns into a coordinate system that is more centered around the fundamental structures found in music. Rather than using the entire frequency spectrum for each piece of music, we can train our classification models using the coordinates given by this system. This will help us develop our classifiers by providing a low-dimensional space for them to train in.

### 4.2.1 Band Classification

In this test case, we tried to build a model to differentiate the music of three bands from different genres: AC/DC, Frank Sinatra, and Beehtoven. The goal was to determine which band each piece came from. Running our model, we saw that SVM was able to separate the training data fairly well. After generating the model, only around 10% of the training data was misclassified. However, we can see that this result does not generalize well as the model's accuracy on the testing data was only around 45%, slightly above random chance. DLT performed similar to SVM with a low training error rate, only 5%, but a low accuracy on the testing data, only 33%, suggesting this model also does not generalize well. On the other hand, KNN did a much better job at segregating the training data, developing a model which classifies every piece of data in the training set correctly. Futhermore, KNN had significantly better performance on the testing data set than either SVM or DLT, giving around 65 % accuracy. As we can see in **Figure 4**, SVM was fairly consistent, correctly classifying and misclassifying each class of music about the same percentage of the time, while KNN got a bulk of its correct classifications from class two (Beethoven). Interestingly, Decision Tree classifications almost never predicted case 2 (Beethoven) and had a large preference for predicting case 3 (Sinatra).

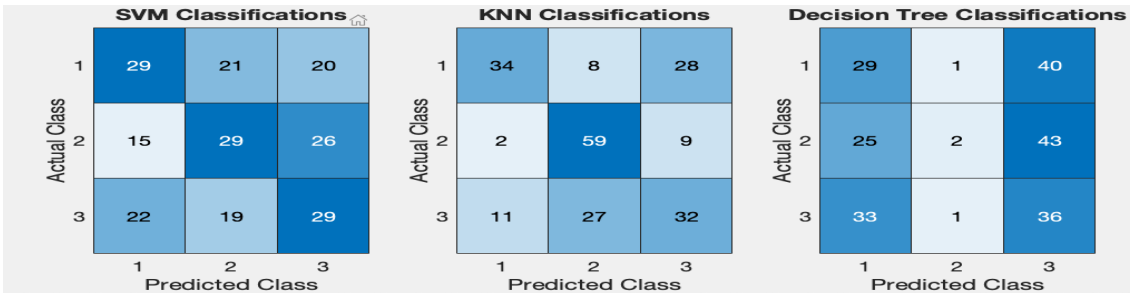


Figure 4: Case 1 classifications for each of the supervised learning algorithms. Entry  $(i, j)$  of each figure counts the number of test cases that were classified as band  $i$  and were actually band  $j$ . In this case, class one is AC/DC, class two is Beethoven, and class 3 is Sinatra.

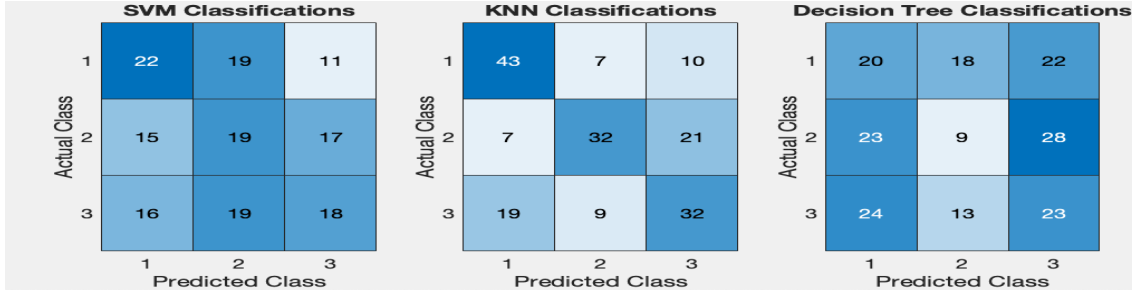


Figure 5: Case 2 classifications for each of the supervised learning algorithms. Entry  $(i, j)$  of each figure counts the number of test cases that were classified as band  $i$  and were actually band  $j$ . In this case, class one is Alice in Chains, class two is Pearl Jam, and class 3 is Sound Garden.

#### 4.2.2 The Case For Seattle

In this test case, we tried to build a model to differentiate the music of three bands from the same genre, grunge/rock, and all from the Greater Seattle area: Alice in Chains, Pearl Jam, and Soundgarden. The goal was to determine which band each piece came from. Running our model, we saw that SVM was not able to separate the training data fairly well. After generating the model, over 35% of the training data was misclassified. However, we can see that unlike in the previous case, this result does generalize well: the model's accuracy is about 62 %, approximately the same as what it was on the training data. DLT performed better than SVM on the training data, only misclassifying 5% of the sample, but did roughly the same on the testing set with an accuracy of around 59%. Once again, KNN did a much better job at segregating the training data, developing a model which classifies every piece of data in the training set correctly. Furthermore, KNN had significantly better performance on the testing data set than either SVM or DLT, giving around 77 % accuracy. From **Figure 5**, we can see that once again SVM is much more consistent than the other two methods, correctly predicting and misclassifying each of the possible classes at the same rate. Furthermore, once again *KNN* is better at predicting some classes than the others, almost never misclassifying a piece in class one (Alice in Chains). In this test case, Decision Tree classification was much more consistent, although it seems to have particular difficulty classifying pieces from class two (Pearl Jam). Interestingly, the models performed better in this case than any of the others.

#### 4.2.3 Genre Classification

In this test case, we tried to build a model to differentiate the music of a variety of bands from the same genre. The goal was to determine which genre each piece of music fell in to: classical (Bach, Mozart, Beethoven, Vivaldi), heavy metal (Metallica, Iron Maiden, Black Sabbath), or country (Johnny Cash, Zac Brown Band, and Merle Haggard). Running our model, we saw that SVM was able to separate the training data fairly well. After generating the model, only around 13% of the training data was misclassified. However, we can see that this result does not generalize well as the model's accuracy on the testing data was only around 43%, slightly above random chance. DLT performed similar to SVM with a low training error rate, only 5%, but a low accuracy on the testing data, only 34%, suggesting this model also does not generalize well. On the other hand, KNN did

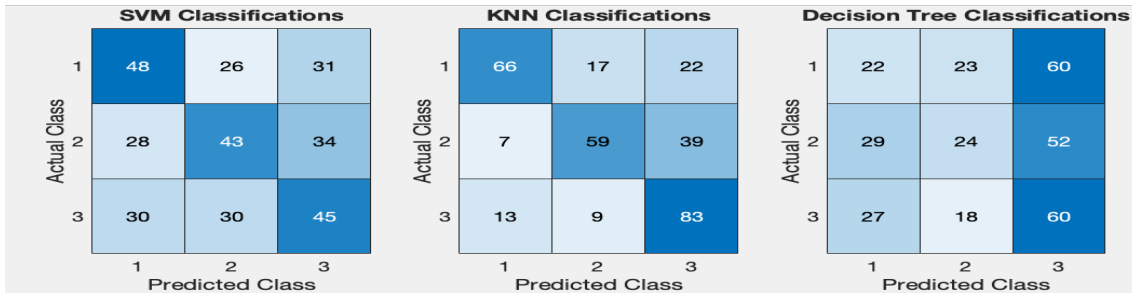


Figure 6: Case 3 classifications for each of the supervised learning algorithms. Entry  $(i, j)$  of each figure counts the number of test cases that were classified as band  $i$  and were actually band  $j$ . In this case, class one is classical music, class two is country music, and class 3 is heavy metal music.

a much better job at segregating the training data, developing a model which classifies every piece of data in the training set correctly. Furthermore, KNN had significantly better performance on the testing data set than either SVM or DLT, giving around 66 % accuracy. As we can see in **Figure 6**, SVM was fairly consistent, correctly classifying and misclassifying each class of music about the same percentage of the time, while KNN got a bulk of its correct classifications from class three (heavy metal). Interestingly, the Decision Tree Classification heavily favored the prediction of class three (heavy metal) regardless of the actual class of the data. This is likely what led to its poor accuracy.

## 5 Summary and Conclusions

From the Yale Faces example, we saw that the SVD cannot handle translated data very well. Without the preprocessing steps required to carefully center the face, the dimension of the uncropped images under SVD quickly blew up. Not only did this make it difficult to interpret the underlying features being captured by each mode, but it drastically increased the computational work needed to manipulate this system. As we saw in the Music Classification section, if we try to apply this SVD decomposed system to a classification problem, we often run into slow runtimes or overfitting. Instead, if we invest a bit of time into preprocessing our data, we can generate a low-rank approximation. From this approximation, we can pick out only the principal coordinates that best separate the data. Rather than working with the entire dataset, we only have to work with the coordinates of our data in this reduced-dimension space. Not only will this reduce the likelihood of overfitting, but it also makes our system easier to work with computationally. This may allow us to train on larger datasets and improve the quality of our predictions.

## 6 Appendix A

Below is a brief summary of the MATLAB functions I used during this project and their functions.

**svd(X)**: Calculates the Singular Value Decomposition of the given matrix  $X$ , returning  $U, \Sigma, V$ .

**fitcecoc(training data, labels)**: Generates a predictive model using multi-class SVM, training it with the given data and its associated labels.

**fitcknn(training data, labels)**: Generates a predictive model using K-Nearest Neighbors, training it with the given data and its associated labels.

**fitctree(training data, labels)**: Generates a predictive model using Decision Tree Learning, training it with the given data and its associated labels.

**fft(x)**: Computes the frequency content of the given vector.

**pcolor(X)**: plots a colored representation of the given matrix.

**flipud(X)**: Flips a matrix top to bottom to be compatible with pcolor.



## 7 Appendix B

### 7.1 Main Methods

#### 7.1.1 Part 1: Yale Faces

```
1 % AMATH 482 HW 4
2
3 % Zachary McNulty
4
5 % NOTES: benefits of alignment is that we can drastically reduce the
6 % number of modes
7 % we have. SVD SUCKS with translated data so the uncropped images raise
8 % the supposed
9 % rank of the system because SVD cannot recognize the translated
10 % images as
11 % the same. Takes time for preprocessing.
12
13
14 % Convert of music into frequency space to avoid the issues with
15 % translation that interferes
16 % With SVD
17
18 %% Part 1: Yale Faces CROPPED
19
20 clear all; close all; clc;
21
22 % path to hw folder
23 % system('cd ~/Desktop/AMATH482/hw/hw4')
24
25 % cropped vs uncropped shows how poor SVD does with translated images
26 % and
27 % the importance of pre-processing.
28 folders = dir('~/Desktop/AMATH482/hw/hw4/input_files/CroppedYale/yale*');
29
30 % MY FACES!
31 F_ave = zeros(32256, 38);
32 F = zeros(32256, 2415);
33
34 index = 1;
35 index2 = 1;
36 for j = 1:length(folders)
37     path = strcat(folders(j).folder , '/', folders(j).name, '/*.pgm');
38     files = dir(path);
39
40     for k = 1:length(files)
41         face = imread(strcat(files(k).folder , '/', files(k).name));
42
43         % average all the images for each person
44         % Since all the images have the same positional location of the
45         % face this seems to help balance out the lighting in the
46         % photograph.
47         F(:, index2) = reshape(double(face), size(face, 1)*size(face, 2),
48             , 1);
49         index2 = index2 + 1;
50
51         if k == 1
52             ave_face = double(face);
53         else
54             ave_face = ave_face + double(face);
55         end
56     end
57 end
```

```

49         end
50     end
51     ave_face = ave_face ./ length(files);
52
53     % each row is a measurement of a given pixel across many trials (i.
54     % e.
55     % photos of people)
56     F_ave(:, index) = reshape(ave_face, size(ave_face, 1)*size(ave_face
57     , 2), 1);
58     index = index + 1;
59 end
60
61 % Demean the data
62 mean_ave = mean(F_ave, 2);
63 F_ave = F_ave - mean_ave;
64 mean_F = mean(F, 2);
65 F = F - mean_F;
66
67 imheight = size(face, 1);
68 imwidth = size(face, 2);
69
70 %% Calculate the Singular Value Decomposition
71
72 % [u, s, v] = svd(F_ave, 'econ');
73 % [u2, s2, v2] = svd(F, 'econ');
74
75 %% Interpretation of U, Sigma, V
76
77 % set which group of images we are working with, averaged or not
78 U = u;
79 S = s;
80 V = v;
81 X = F_ave;
82 m = mean_ave;
83
84 %%
85 % The columns of U represent an orthogonal basis for the codomain of
86 % our data matrix X.
87 % Not just any basis, however, the BEST basis for the space where each
88 % successive column captures as much of the variance in the data as
89 % possible.
90 % In this case, as the codomain of our matrix is the space in which
91 % column of
92 % the matrix exists, as each column is a face in our data set the
93 % codomain is
94 % "face space" Thus, it follows that
95 % U is the "optimal" basis for face-space where the SVD's idea of what
96 % the important features
97 % that make up a face are will be those that capture the most variance.
98 % Furthermore, the columns of V then store the
99 % coordinates of each of our data measurements (i.e. each face) within
100 % this
101 % face space. What linear combination of these principal components
102 % (outputted by the code below).
103 num_faces = 5;
104 for col = 1:num_faces
105     figure(2)
106
107     subplot(1, num_faces, col)
108     next_face = reshape(U(:, col), [imheight, imwidth]);

```

```

103     pcolor(flipud(next_face)), shading interp, colormap(gray);
104     title(strcat("Principal Component ", num2str(col)));
105     set(gca, 'xticklabel', [])
106     set(gca, 'yticklabel', [])
107     set(gca, 'fontsize', 15)
108 end
109
110
111
112 %% Plot singular Values
113
114 figure(4)
115 plot(diag(S) ./ max(diag(S)), 'r.', 'markersize', 20)
116 %xticks(1:round(length(S)/10):length(S))
117 xticks(0:20:240)
118 xlim([0, 250])
119 ylim([0,1])
120 title('Normalized Singular Values')
121 ylabel('\sigma_j')
122 xlabel('index (j)')
123 set(gca, 'fontsize', 20)
124
125 %% Choosing rank r
126
127 % We will choose the r that captures at least 90% of the energy of the
128 % system
129
130 energy = 0;
131 total = sum(diag(S));
132 % how much energy we want our modes to capture.
133 % 75% does alright; 90% does very good.
134 threshold = 0.30;
135 r = 0;
136 while energy < threshold
137     r = r + 1;
138     energy = energy + S(r,r)/total;
139 end
140
141
142 %% rank r approximation of face space
143
144 % low-rank approximation of our faces stored in F.
145 X_r = U(:, 1:rank) * S(1:rank, 1:rank) * (V(:, 1:rank)');
146
147 faces_to_plot = 10;
148
149 for j = 1:faces_to_plot
150     figure(5)
151     subplot(211)
152     imshow(uint8(reshape(X_r(:, j) + m, imheight, imwidth)));
153     title(strcat("Rank ", num2str(rank), " approximation"))
154     set(gca, 'fontsize', 15);
155     subplot(212)
156     imshow(uint8(reshape(X(:, j) + m, imheight, imwidth)));
157     title('Original Image')
158     set(gca, 'fontsize', 15);
159 end
160
161
162 %% Part 1: Yale Faces UNCROPPED
163

```

```

164 clear all; close all; clc;
165
166 % NOTE: in this case averaging might not work as well as the previous
      case
167 % because the photos are no longer aligned.
168
169 % cropped vs uncropped shows how poor SVD does with translated images
      and
170 % the importance of pre-processing.
171 folders = dir('~/Desktop/AMATH482/hw/hw4/input_files/
      yalefaces_uncropped/subject*');
172
173 % MY FACES!
174 F_ave = zeros(77760, 15);
175 F = zeros(77760, 165);
176
177 index = 1;
178 index2 = 1;
179 for j = 1:length(folders)
180     path = strcat(folders(j).folder , '/', folders(j).name, '/subject*'
      );
181     files = dir(path);
182
183     for k = 1:length(files)
184         face = imread(strcat(files(k).folder , '/', files(k).name));
185
186         % average all the images for each person
187         % Since all the images have the same positional location of the
188         % face this seems to help balance out the lighting in the
189         % photograph.
190         F(:, index2) = reshape(double(face), size(face, 1)*size(face, 2)
      , 1);
191         index2 = index2 + 1;
192
193         if k == 1
194             ave_face = double(face);
195         else
196             ave_face = ave_face + double(face);
197         end
198     end
199     ave_face = ave_face ./ length(files);
200
201     % each row is a measurement of a given pixel across many trials (i.
      e.
202     % photos of people)
203     F_ave(:, index) = reshape(ave_face , size(ave_face, 1)*size(ave_face
      , 2), 1);
204     index = index + 1;
205 end
206
207 % Demean data?
208 mean_ave = mean(F_ave, 2);
209 F_ave = F_ave - mean_ave;
210 mean_f = mean(F, 2);
211 F = F - mean_f;
212
213 imheight = size(face, 1);
214 imwidth = size(face, 2);
215
216 %% Calculate the Singular Value Decomposition
217

```

```

218 [u, s, v] = svd(F_ave, 'econ');
219 [u2, s2, v2] = svd(F, 'econ');
220
221
222 %% Interpretation of U, Sigma, V
223
224 % set which group of images we are working with, averaged or not
225 U = u2;
226 S = s2;
227 V = v2;
228 X = F;
229 m = mean_f;
230
231
232 %%
233 % The columns of U represent an orthogonal basis for the codomain of
    our data matrix X.
234 % Not just any basis, however, the BEST basis for the space where each
235 % successive column captures as much of the variance in the data as
236 % possible.
237 % In this case, as the codomain of our matrix is the space in which
    column of
238 % the matrix exists, as each column is a face in our data set the
    codomain is
239 % "face space" Thus, it follows that
240 % U is the "optimal" basis for face-space where the SVD's idea of what
    the important features
241 % that make up a face will be those that capture the most variance.
242 % Furthermore, the columns of V then store the
243 % coordinates of each of our data measurements (i.e. each face) within
    this
244 % face space. What linear combination of these principal components
245 % (outputted by the code below).
246
247 num_faces = 5;
248 for col = 1:num_faces
249     figure(2)
250
251     subplot(1, num_faces, col)
252     next_face = reshape(U(:, col), [imheight, imwidth]);
253     pcolor(flipud(next_face)), shading interp, colormap(gray);
254     title(strcat("Principal Component ", num2str(col)));
255     set(gca, 'xticklabel', [])
256     set(gca, 'yticklabel', [])
257     set(gca, 'fontsize', 15)
258 end
259
260
261
262
263 %% Plot singular Values
264
265 figure(4)
266 plot(diag(S) ./ max(diag(S)), 'r.', 'markersize', 20);
267 ylim([0,1])
268 title('Normalized Singular Values')
269 ylabel('\sigma_j')
270 xlabel('index (j)')
271 set(gca, 'fontsize', 20);
272 %% Choosing rank r
273

```

```

274 % We will choose the r that captures at least 90% of the energy of the
275 % system
276
277 energy = 0;
278 total = sum(diag(S));
279 % how much energy we want our modes to capture.
280 % 75% does alright; 90% does very good.
281 threshold = 0.9;
282 r = 0;
283 while energy < threshold
284     r = r + 1;
285     energy = energy + S(r,r)/total;
286 end
287
288
289 %% rank r approximation of face space
290
291 rank = r;
292
293 % low-rank approximation of our faces stored in F.
294 X_r = U(:, 1:rank) * S(1:rank, 1:rank) * (V(:, 1:rank)') + m;
295
296 faces_to_plot = 10;
297
298 for j = 1:faces_to_plot
299     figure (5)
300     subplot(211)
301     imshow(uint8(reshape(X_r(:, j), imheight, imwidth)));
302     title('Low rank approximation')
303     subplot(212)
304     imshow(uint8(reshape(X(:, j) + m, imheight, imwidth)));
305     title("Original Image")
306     pause(1);
307 end

```

### 7.1.2 Part 2: Music Classification

```
1 %% Part 2: Music Classification
2
3
4 %% Test case 1: Different Bands from Different Genres
5
6
7 % Collect Data
8 clear all; close all; clc;
9
10 % path to hw folder;
11 % set the case you want to work on here: part1,part2,part3
12 folders = dir('~/Desktop/AMATH482/hw/hw4/input_files/music_files/part3
    /*');
13
14 % number of 5 second long samples to be taken from each song
15 % NOTE: takes actually twice this number of samples as half the samples
    are
16 % placed in the training set and the other half in the validation set.
17 samples_per_song = 15;
18 num_songs = 21; % total number of songs across all groups
19 sample_length = 5; % in seconds
20 skip_period = 15; % skip the first "skip period" seconds of song to
    avoid aampling silence at beginning of song.
21 %sample_rate = 2; % take every other "sample_rate"th point.
22
23 % Fs is always 44100 across all the songs we downloaded based on
24 % the download procedure
25 Fs = 44100;
26
27
28 % Each row is a sample
29 % each data sample will be the frequency content of a given song
30 training_data = zeros(Fs*sample_length + 1, samples_per_song *
    num_songs);
31 %training_data = zeros(262152, samples_per_song * num_songs); %
    spectrogram
32 training_labels = zeros(samples_per_song * num_songs, 1);
33
34 validation_data = zeros(Fs*sample_length + 1, samples_per_song *
    num_songs);
35 %validation_data = zeros(262152, samples_per_song * num_songs); %
    spectro
36 validation_labels = zeros(samples_per_song * num_songs, 1);
37
38 index = 1;
39 group_num = 1;
40
41 tic
42 for j = 1:length(folders)
43
44     % skip all hidden folders within the directory
45     if startsWith(folders(j).name, '.')
46         continue
47     end
48
49     path = strcat(folders(j).folder , '/', folders(j).name, '/*.mp3');
50     files = dir(path);
51
52     for k = 1:length(files)
```

```

53     song_path = strcat(files(k).folder, "/", files(k).name);
54
55
56     start = skip_period * Fs; % skip first "skip_period" seconds of
        song.
57     finish = inf; % sample until the end of audio file
58
59     % Fs is the sampling rate and Y is the amplitude at each point
        in
60     % the recording
61     [Y, Fs] = audioread(song_path, [start, finish]);
62
63     % Y is a stereo measurement (includes measurements for both
        left
64     % and right speakers) so we average these two to get a single
65     % measurement
66
67     Y = mean(Y, 2).';
68     %p8 = audioplayer(Y,Fs); playblocking(p8);
69
70
71     % randomly make training data set and cross validation data set
72     for s = 1:samples_per_song
73
74         % randomly choose a sample starting point
75         sample_start = randi(length(Y) - sample_length*Fs, [2,1]);
76
77         % from the given starting point, take a sample of
78         % 'sample_length' seconds and store the corresponding label
79         % Take the fft of the sample to convert to frequency space
80
81         training_data(:, index) = fft(Y(sample_start(1):
            sample_start(1) + sample_length * Fs)).';
82         %spec = spectrogram(Y(sample_start(1):sample_start(1) +
            sample_length * Fs));
83         %training_data(:, index) = reshape(spec, size(spec,1) *
            size(spec,2), 1);
84         training_labels(index) = group_num; % i.e. this song came
            from band j
85
86         % take another sample to be used for validation
87         validation_data(:, index) = fft(Y(sample_start(2):
            sample_start(2) + sample_length * Fs)).';
88
89         %spec = spectrogram(Y(sample_start(2):sample_start(2) +
            sample_length * Fs));
90         %validation_data(:, index) = reshape(spec, size(spec,1) *
            size(spec,2), 1);
91         validation_labels(index) = group_num;
92         index = index + 1;
93
94         %         if s == 1 && k == 1
95         %             figure(10)
96         %             t = 0:1/Fs:5;
97         %             subplot(3, 1, group_num);
98         %             plot(t, Y(sample_start(1):sample_start(1) +
sample_length * Fs), 'r');
99
100        %             xlim([0,5])
101        %
102        %             if group_num == 1

```



```

103 %             title('AC DC')
104 %         elseif group_num ==2
105 %             title('Beethoven')
106 %         elseif group_num == 3
107 %             xlabel('Time (s)')
108 %             title('Frank Sinatra')
109 %         end
110 %         ylabel('Amplitude')
111 %         set(gca, 'fontsize', 15);
112 %     end
113
114     end
115
116     end
117
118     group_num = group_num + 1;
119
120 end
121 toc
122
123
124 %% Find Modes which are best at separating groups.
125
126 % Extract the individual groups (i.e. bands in this case) from the
    training data
127 group1 = training_data(:, training_labels == 1);
128 group2 = training_data(:, training_labels == 2);
129 group3 = training_data(:, training_labels == 3);
130
131 % SVD the entire dataset to find principal components
132 % of "music space"
133 mean_td = mean(training_data, 1);
134 [u, s, v] = svd(training_data - mean_td, 'econ'); % mean subtract
135
136 % Plot singular values: which are relevant?
137
138 singular_values = diag(s) / max(diag(s));
139 plot(singular_values, 'r.', 'markersize', 20)
140 title('Normalized Singular Values')
141 ylabel('\sigma_j')
142 xlabel('index j')
143
144 % Find the modes that best separate the data by projecting each of our
145 % groups onto the individual components. Note that each column of V
    gives
146 % the coordinates of each data measurement (row in X = training_data)
    onto
147 % the corresponding principal component (column in U). So entry (i,j)
    of V
148 % gives the weighting/"importance" of principal component
149
150
151 %% Low rank approximation
152
153 % Calculate the low-rank approximation that captures at least 90% of
    the
154 % energy of the system.
155 energy_threshold = 0.90;
156 rank = find(cumsum(singular_values / sum(singular_values)) >=
    energy_threshold, 1);
157

```

```

158 X_r = u(:, 1:rank) * s(1:rank, 1:rank) * (v(:, 1:rank)') + mean_td;
159
160 %% Test Quality of reconstruction
161
162 % test = real(iff t(X_r(:, 1))); % all complex parts are zero and
    interfere with playing music
163 % p8 = audioplayer(test, Fs);
164 % playblocking(p8);
165
166 %% Train a SVM model
167 clc;
168
169 % SVM Cannot handle complex numbers so we simply take the absolute
    value
170 % To convert the training data to amplitude (all real).
171
172 % In this case, as our data measurements are stored as columns, the
    columns of
173 % U store the fundamental structure/modes of our system (they form a
    basis
174 % for our domain : frequency / music space) while the columns of SV'
    give the
175 % coordinates of each data measurement (i.e. each song) within this
176 % frequency space. Thus, these coordinates define the underlying
    structure
177 % of our system and we can use them to classify.
178
179 % Specifically, we will use V* as our "coordinates" and project future
180 % data measurements onto the columns of U*S
181
182 vstar = v';
183 xtrain = [real(vstar(1:rank, :)) ; imag(vstar(1:rank, :))].';
184
185 %SVM model
186 Mdl_svm = fitcecoc(xtrain, training_labels);
187
188 % K nearest neighbors model
189 Mdl_knn = fitcknn(xtrain, training_labels);
190
191 % Naive Bayes Model
192 %Mdl_nb = fitcnb(xtrain, training_labels);
193
194 % Decision Tree Classification model
195 Mdl_tree = fitctree(xtrain, training_labels);
196
197
198 % Computes the classification Error rate in classification for our
    training
199 % data set using the classification conditions defined by our model.
200 error_rate_training_svm = resubLoss(Mdl_svm)
201 error_rate_training_knn = resubLoss(Mdl_knn)
202 %error_rate_training_nb = resubLoss(Mdl_nb)
203 error_rate_training_tree = resubLoss(Mdl_tree)
204
205 %% Test how good the model is at predicting labels
206
207
208 % since our model was trained on the coordinates given by V which are
    the coordinates
209 % of our data with respect to the basis US (for frequency space), we
    have to

```

```

210 % to get the coordinates of our validation data with respect to this
211 % basis:  $US * coordinates\_in\_US = data$ . Since  $U$  is unitary, this yields
212 % simply:  $coordinates\_in\_US = S^{-1} U' * data$ 
213 Sinv = diag(1 / diag(s));
214 coordinates_in_US = (u')*validation_data;
215 xtest = [real(coordinates_in_US(1:rank, :)) ; imag(coordinates_in_US(1:
    rank, :))].';
216
217 % we only used the first  $r = rank$  coordinates to determine our
218 % classification so we again do that.
219 predictions_svm = predict(Mdl_svm, xtest);
220 predictions_knn = predict(Mdl_knn, xtest);
221 %predictions_nb = predict(Mdl_nb, xtest);
222 predictions_tree = predict(Mdl_tree, xtest);
223 accuracy_svm = sum(predictions_svm == validation_labels) / length(
    validation_labels)
224 accuracy_knn = sum(predictions_knn == validation_labels) / length(
    validation_labels)
225 %accuracy_nb = sum(predictions_nb == validation_labels) / length(
    validation_labels)
226 accuracy_tree = sum(predictions_tree == validation_labels) / length(
    validation_labels)
227
228
229
230
231 %% Find which kinds of misclassifications were common.
232
233 figure(10)
234 subplot(131)
235 classification_heatmap(validation_labels, predictions_svm);
236 title('SVM Classifications')
237 subplot(132)
238 classification_heatmap(validation_labels, predictions_knn);
239 title('KNN Classifications')
240 subplot(133)
241 classification_heatmap(validation_labels, predictions_tree);
242 title('Decision Tree Classifications')

```

## 7.2 Helper Methods

```
1 function x = classification_heatmap(actual, predictions)
2 % returns a heatmap where columns are "classified as" and rows are "
  actual"
3 x = zeros(3,3);
4 for j = 1:3
5     for k = 1:3
6         x(j, k) = sum((actual == j) .* (predictions == k));
7     end
8 end
9
10 heatmap(x);
11 ylabel('Actual Class')
12 xlabel('Predicted Class')
13 set(gca, 'fontsize', 20);
14 colorbar(gca, 'off')
15
16 end
```