

Machine Learning for Music Classification: A Test of Linear Discriminant Analysis and k-Nearest Neighbors in Three Iterations

Dana Korssjoen

3/6/2020

Abstract

This analysis applies a machine learning approach to the task of distinguishing between musical artists and genres using only computational tools. The model developed in this report uses singular value decomposition (SVD), linear discriminant analysis (PCA), and a clustering algorithm (k-nearest neighbors) to accomplish the task. The model is then tested against 3 distinct datasets, each with different structures, achieving a relatively high degree of accuracy in the best case (80%) and a decent degree of accuracy in the worst case (67%).

1 Introduction and Overview

One of the most interesting applications of math is trying to get computers to do things that humans can already do quite easily. I'm speaking, of course, about machine learning.

In this report, the “thing” that we want to teach the computer to do is to recognize music. The mathematics of the model are discussed at length in later sections, but they rely fundamentally on singular value decomposition (SVD), linear discriminant analysis (PCA), and a clustering algorithm (k-nearest neighbors). Before any learning can be done, however, an algorithm to generate the dataset of songs is developed, and these songs are transformed into their spectrograms to facilitate computation.

The particular question asked by this report is, how do different datasets respond to the algorithm developed here? Each dataset has three classes of music, and we break the datasets into 3 cases: case I, where each of the three classes contains songs from different musical artists in different genres; case II, where all three classes come from the same genre, but different artists; and case III, where the three classes come from different genres, and each class contains songs from multiple artists. In each case, we test whether the model can correctly classify samples from songs that were not included in the training data.

All of the music used in this analysis was downloaded for free from the Free Music Archive. Case I features songs from The Tudor Consort, a classical choir ensemble; Derek Clegg, a folk/acoustic artist; and King Imagine, an artist who makes electronic/ambient music. Again, all of these artists come from different genres. Case II features three artists, all from the same genre: King Imagine, as previously discussed; Livio Amato, who produces ambient/soundscape music that combines classical piano with electronic elements; and Soularflair, who produces deep electronic music. Each of these artists occupies the same genre, but there are noticeable differences in their instrumentation, rhythms, and tones. I recommend all three artists if you ever need some good studying music. Finally, case III features multiple artists within three different genres. The first genre is classical choral music; then Balkan big band (which is arguably the best genre of music in existence); and then finally, ambient/soundscape/-electronic, featuring the artists from case II. I am deeply grateful to all of these artists for making their music publicly available.

2 Theoretical Background

2.1 Singular Value Decomposition

A necessary step in this analysis is transforming the data under SVD. Because this method has been discussed at length in other reports, I won't go into excessive detail; it suffices to mention that SVD can be used to extract the key components \mathbf{U} , $\mathbf{\Sigma}$, \mathbf{V} , and that we can capture the majority of the energy of the data using only the first several principal components. It is this projected, reduced dataset that we will perform LDA on in the next section. We can see the share of total energy captured by each mode for the SVDs of datasets I and II in figure 3.

2.2 Linear Discriminant Analysis

Linear Discriminant Analysis (LDA) is a method of projecting data onto new a basis that maximizes the distance between the inter-class data while minimizing the distance between intra-class data. In simpler terms, LDA is intended to create distinct clusters of same-classed data.

Mathematically, the ideal projection \mathbf{w} will take the form

$$\mathbf{w} = \arg \max_{\mathbf{w}} \frac{\mathbf{w}^T \mathbf{S}_b \mathbf{w}}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}} \quad (1)$$

where \mathbf{S}_b and \mathbf{S}_w are the between-class and within-class variances, respectively, computed by

$$\mathbf{S}_b = \sum_{i=1}^C (\mu_i - \mu)(\mu_i - \mu)^T \quad (2)$$

and

$$\mathbf{S}_w = \sum_{i=1}^C \sum_{\mathbf{x}} (\mathbf{x} - \mu_i)(\mathbf{x} - \mu_i)^T \quad (3)$$

where C is the number of classes, μ is the average value across the whole dataset, and μ_i is the average value across class i .

Once these quantities are computed, we can find the basis \mathbf{w} that achieves the desired result by solving the generalized eigenvalue problem

$$\mathbf{S}_b \mathbf{w} = \lambda \mathbf{S}_w \mathbf{w} \quad (4)$$

where the maximum eigenvalue(s) λ and the associated eigenvector(s) gives the desired projection basis. Later sections of this paper will give concrete examples of data being projected onto this basis.

3 Algorithm Implementation and Development

All line numbers given in this section refer to the line numbers of the MATLAB code, given in section 7, “Appendix B: MATLAB Code.”

3.1 Generating the Dataset

Lines 181 - 256 generate the dataset. Not shown here (though uploaded to Github) is the set of very large matrices of urls that are used to download each individual song in each dataset. This matrix is loaded on line 186, then the minimum sampling rate for the songs in the dataset is computed by `getMinFs()`, which is a function that simply loops over each element of the dataset. That sampling rate is then used to downsample other songs as needed (to standardize the dataset), at which point, songs are cut into 5 second samples. Each song generates 3 samples, taken from 3 different points in the song, to increase the size of the dataset and accurately represent all features of that song.

All three datasets, as well as their labels, are then saved to be used in analysis.

Once the datasets are loaded into the main script on line 2, they are immediately split into testing and training sets with a ratio of 80/20 by the function given on lines 25-34. As seen in the results section, splitting the dataset at a specific index, rather than randomly choosing data points to be in the training set, may have introduced bias. This method was chosen consciously, however, because if data points were chosen randomly, there is extremely high probability that most of the data points would be from songs that the classifier had trained on, which would have introduced even more bias.

3.2 Computing the Spectrograms

Line 11 refers to a function on lines 75-92 that generate the matrix of spectrograms for the entire dataset. This is done using MATLAB’s built-in spectrogram function, then unrolling each spectrogram into a column vector, so that each column of the resulting matrix is an observation from the dataset. An example signal and its spectrogram are given in figures 1 and 2.

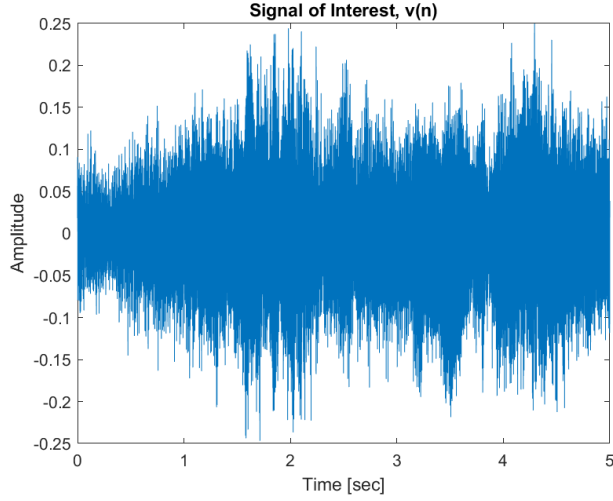


Figure 1: Frequency over time of a single data point.

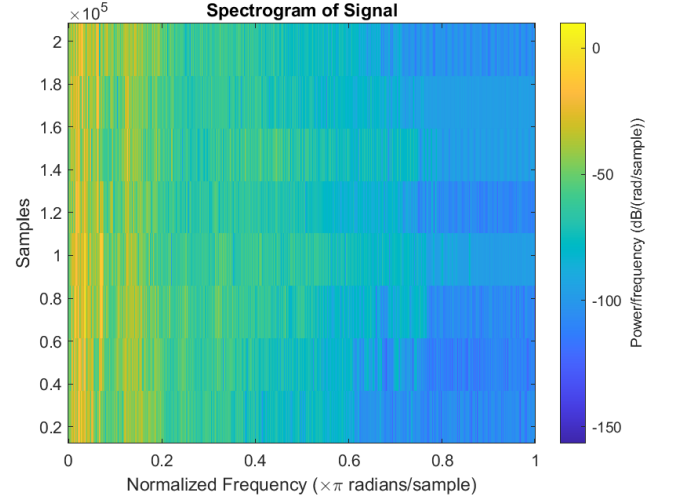


Figure 2: The spectrogram of that same data point.

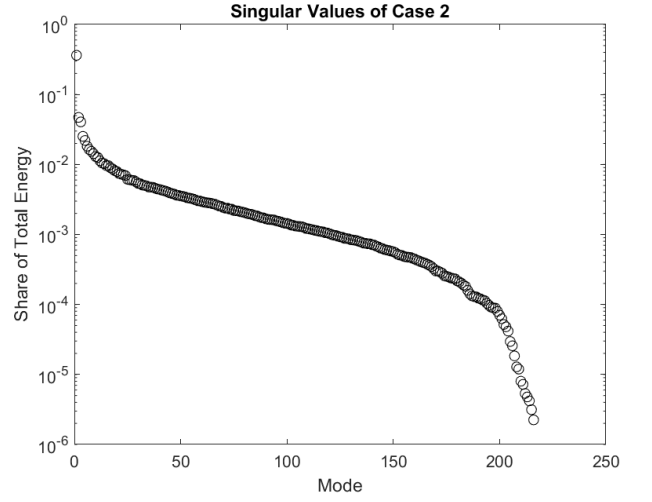
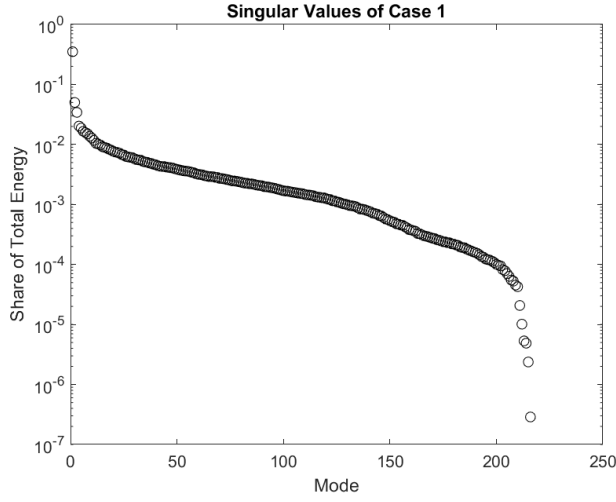


Figure 3: Energy of the singular values of case I (left) data and case II (right) data. Energy is depicted as the share of total energy represented by that mode.

3.3 Linear Discriminant Analysis

Once the spectrograms have been computed, it is time to perform LDA on line 13, which refers to a function on lines 36-73. In order to do this, the number of features is chosen based on the plots in figure 3. The number 25 was chosen because the distribution of energy among the modes is heavy-tailed, but higher numbers were found experimentally to capture too much noise and thus increase error. After the projection onto the first 25 principal components is computed, the within- and between-class variances are found. As discussed in the theoretical section, we use these to compute the best bases for our projection, hoping to maximize the distance between the means of clusters while minimizing within-class variance. The projection of the data onto these bases is returned. The results are graphed in Section 4, “Computational Results.” The choice to project the data into 2-dimensional, rather than 1-dimensional space was made because there is significant overlap between classes in any 1-dimensional projection.

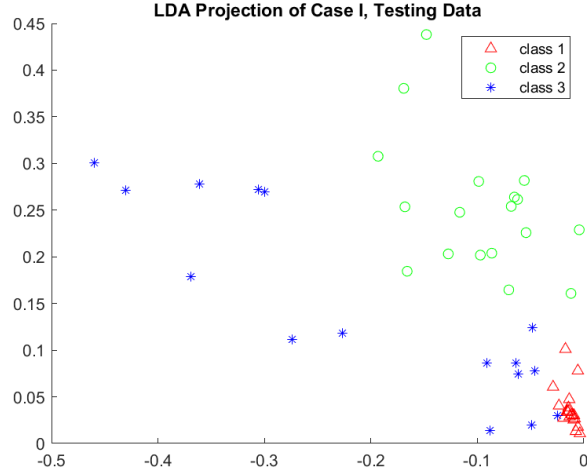


Figure 4: Testing data from case I

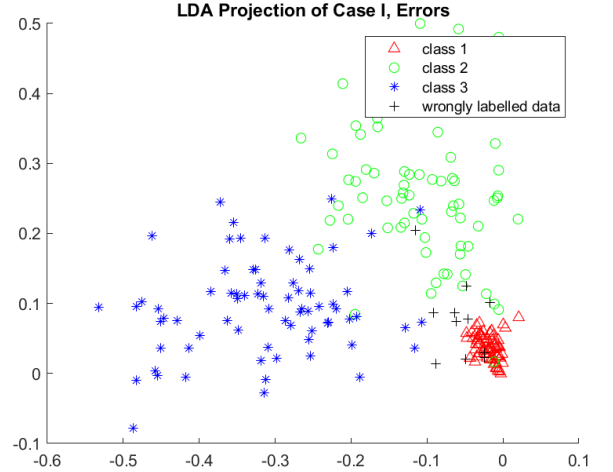


Figure 5: Mislabelled testing data from case I

Predicted	2	1	1	1	1	2	1	2	1	2	1
Actual	1	3	3	3	3	3	3	3	3	3	3

Table 1: Mislabelled points in case I dataset.

3.4 Classification

Finally, classification of the projected test data is performed using the popular k-nearest neighbors (KNN) algorithm, which is called on line 18 and written on lines 94-106. This algorithm finds the k nearest labelled points to the given unlabelled points, then predicts that the mystery point has the mode (most frequent) of their labels. If there is a tie for which labels appears with the highest frequency, the label which is closest to the point is predicted. The value $k = 5$ was selected experimentally. We can anticipate that KNN is an appropriate classification algorithm for this data set, because our intention with LDA was to group each point into the most homogenous same-classed clusters are possible, and to maximize the distance between clusters. Thus, we would ideally anticipate the nearest neighbors of any given point, would be points in the same class. This idea is confirmed by our relatively low error rate, as discussed in the following section.

4 Computational Results

Figures 4 and 5 depict the LDA projection of the data from case I. In figure 4, only the testing data (with the hidden labels) is projected, and in figure 5, all training data, plus the testing data which was mislabelled by the classifier is projected. In table 1, the predicted vs actual classes of the 11 mislabelled testing points are included. As can be seen in the table, many data points from class 3 were mislabelled. In figures 4 and 5 we can see why this might be the case, as many of the mislabelled class 3 points are confoundingly close to the class 1 cluster. In fact, it appears that the testing points for class 3 are distributed differently than the training points, which suggests there may have been error in the way the dataset was constructed, interference from the relatively small number of data points, or error resulting from the sampling method.

Analogous figures for case II are included as figures 6 and 7, and table 2 depicts the 18 mislabelled testing points. As in case I, the errors in case II seems to come from the fact that the testing data for case II is distributed differently than the training data, particularly for class 1, where we saw the most errors. The reasons for this are discussed above. Intuitively, it makes sense that this case had more error than case I, however, because the degree of difference between bands in the same genre is typically much smaller than that between bands in different genres. Understandably, then, our classifier had more difficulty distinguishing between bands.

The results for case III are included as figures 8 and 9, and table 3 depicts the 12 mislabelled testing points. This case had particular trouble with samples from class 1, which have high variance in the testing data. The same factors as before may have influenced. Despite that, however, we achieved a relatively equal degree of accuracy between cases I and III. Again, this is likely due to that fact that different genres can have vastly different musical

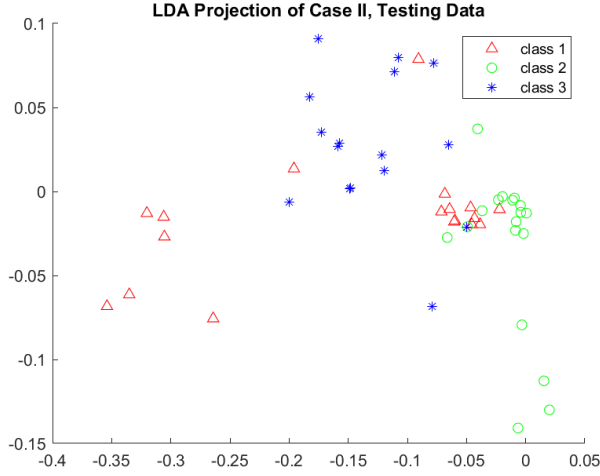


Figure 6: Testing data from case II

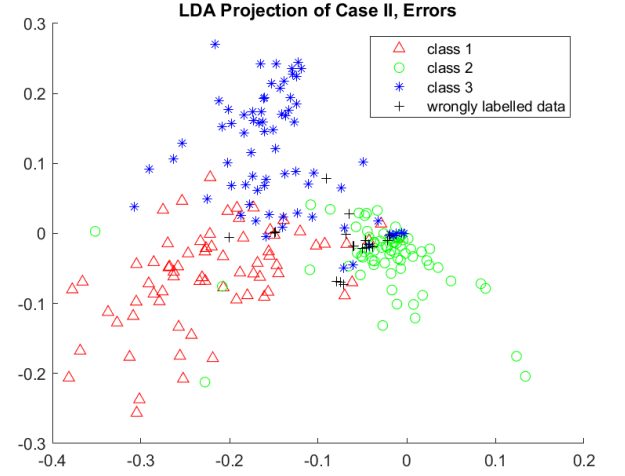


Figure 7: Mislabelled testing data from case II

Predicted	2	2	2	2	2	2	2	2	2	2	2	2	1	1	1	1	1
Actual	1	1	1	1	1	1	1	1	1	2	3	3	3	3	3	3	3

Table 2: Mislabelled points in case II dataset.

characteristics, so the means between the classes would naturally be farther away from one another. Compare this to case II, where same-genre bands are expected to be more similar musically. Note that we cannot conclusively say whether our classifier was more successful with case I or case III, since their total errors was only off by one, which suggests that the more salient feature in our model was which genre a song was from, not which specific musical artist.

All told, the error rate for case I was 11/54, or 20%; for case II was 18/54, or 33%; and for case III was 12/54, or 22%.

5 Summary and Conclusions

As discussed in the computational results section, we may tentatively conclude that the model used in this paper is more equipped to distinguish between genres than between musical artists. There is something that we can conclude confidently, however – the fact that a simple series of mathematical transformations enabled a computer to (somewhat) reliably “recognize” music. The potential applications for this technique, and techniques like it, are profound, particularly if one were to heed the lessons learned in this report by ensuring that datasets are sufficiently large, and training data is obtained in a randomized way in order to minimize error in applications.

6 Appendix A: MATLAB Functions and Implementation

See table 6, which appears underneath the heading for the following section because LaTeX is *very* confusing.

7 Appendix B: MATLAB Code

Predicted	2	2	2	2	2	2	2	2	2	2	1	1	1
Actual	1	1	1	1	1	1	1	3	3	3	3	3	3

Table 3: Mislabelled points in case III dataset.

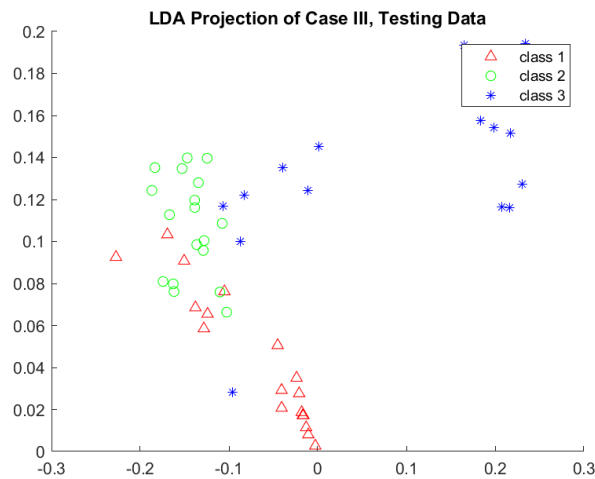


Figure 8: Testing data from case III

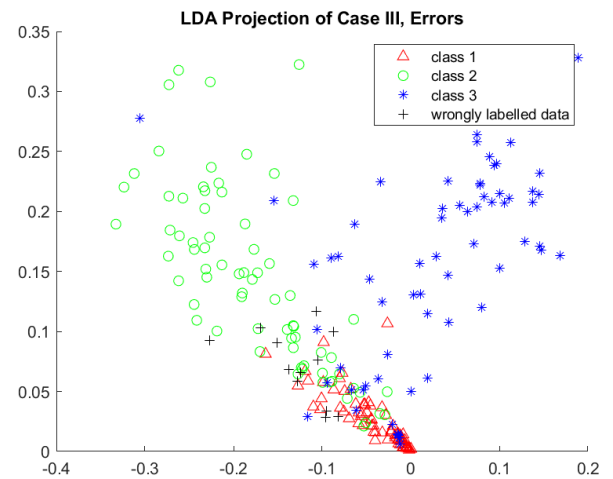


Figure 9: Mislabelled testing data from case III

Function

load(A)
nnz(A)
length(A)
svd(A)
mean(A, dim)
eig(A,B)
maxk(A,k)
size(A,dim)
spectrogram(A)
zeros(n,m)
abs(A)
reshape(A, [n m])
repmat(A, n, m)
mode(A, dim)
webread(A)
downsample(A, n)

Implementation

loads the file or data given by A into the workspace
returns the number of non-zero elements in A
returns the length of vector A
performs a singular value decomposition on A
computes the mean of A along dimension dim
finds the solution to the generalized eigenvalue problem described by A and B
finds the k maximum entries in A
returns the size of A along dimension dim
computes the spectrogram of A; plots it if output is not suppressed by semicolon
returns an nxm matrix of all zeros
returns the absolute value of A
reshapes A into an nxm matrix
replicates matrix A to construct a new matrix of shape nxm
returns the mode of A along dimension dim
reads content from RESTful web service; returns as data
downsamples A by a factor of n

```

1 %% load & split data
2 load('songs.mat')
3 [Atr,Ats,Ltr,Lts] = split_train_test(A,lab);
4
5 %% main body - run this for each dataset
6 X = A; % edit only this line to change dataset
7
8 % split data
9 [Xtr,Xts,Ltr,Lts] = split_train_test(X,lab);
10 % make specs
11 Xsp = spec_vec(Xtr);
12 % make projection
13 [Xtr,U,S,w] = LDA2(Xsp,25,90*.8);
14 % apply to test data
15 Xtssp = spec_vec(Xts);
16 Xts = w'*(U'*Xtssp);
17 % knn
18 plab = knn(Xtr,Ltr,Xts,5);
19 % calculate error
20 errs = nnz(plab~Lts);
21 err_ind = plab~=Lts;
22 err_rate = errs/length(plab);
23
24 %% functions
25 % assumption: obs need to be columns in A
26 function [Atr,Ats,Ltr,Lts] = split_train_test(A,lab)
27     % save 18 out of every 90 for testing (20%)
28     train = 1:72;
29     test = 73:90;
30     Atr = [A(:,train) A(:,train+90) A(:,train+180)];
31     Ltr = [lab(train) lab(train+90) lab(train+180)];
32     Ats = [A(:,test) A(:,test+90) A(:,test+180)];
33     Lts = [lab(test) lab(test+90) lab(test+180)];
34 end
35
36 % LDA2: 2 for 2-dim projection
37 % n is num of each class in the data - classes MUST be consecutive
38 function [result,U,S,w] = LDA2(X, feature,n)
39
40     [U,S,V] = svd(X,'econ');
41     Xp = S*V'; % proj onto principal components
42     U = U(:,1:feature);
43     X1 = Xp(1:feature,1:n);
44     X2 = Xp(1:feature,(n+1):(2*n));
45     X3 = Xp(1:feature,(2*n+1):(3*n));
46
47     % calculate Sw,Sb
48     m1 = mean(X1,2);
49     m2 = mean(X2,2);
50     m3 = mean(X3,2);
51     Sw = 0; % within class
52     for i=1:n
53         Sw = Sw + (X1(:,i)-m1)*(X1(:,i)-m1)';
54         Sw = Sw + (X2(:,i)-m2)*(X2(:,i)-m2)';
55         Sw = Sw + (X3(:,i)-m3)*(X3(:,i)-m3)';
56     end
57     m = (m1+m2+m3)/3;
58     Sb1 = (m1-m)*(m1-m)'; % between class
59     Sb2 = (m2-m)*(m2-m)';
60     Sb3 = (m3-m)*(m3-m)';
61     Sb = Sb1+Sb2+Sb3;
62
63     % LDA
64     [V2,D] = eig(Sb,Sw);
65     [lambda,I] = maxk(abs(diag(D)),2);
66     w1 = V2(:,I(1)); %w1 = w1/norm(w1,2);
67     w2 = V2(:,I(2)); %w2 = w2/norm(w2,2);
68
69     % 2-dim proj
70     w = [w1 w2];
71     v1 = w'*X1; v2 = w'*X2; v3 = w'*X3;

```

```

72     result = [v1,v2,v3];
73 end
74
75 % make specs
76 % X must have observations as cols
77 function [spec] = spec_vec(X)
78     n = size(X,2); % num of observations
79
80     % determine size for memory pre-allocation
81     tmp = spectrogram(X(:,1));
82     sz = size(tmp);
83
84     % make spectrogram matrix
85     spec = zeros(sz(1)*sz(2),n);
86     for j=1:n
87         s = spectrogram(X(:,j));
88         s = abs(s);
89         s = reshape(s, [sz(1)*sz(2),1]);
90         spec(:,j) = s;
91     end
92 end
93
94 % knn
95 function [pred] = knn(tr,lab,ts,k)
96     ntr = size(tr,2);
97     nts = size(ts,2);
98     X = repmat(ts(1,:),ntr,1)-repmat(tr(1,:),1,nts);
99     Y = repmat(ts(2,:),ntr,1)-repmat(tr(2,:),1,nts);
100    l2 = (X.^2 + Y.^2).^(1/2); % Euclid. dist
101
102    % for each col, find k least entries
103    [~,I] = sort(l2,1);
104    pred = lab(I(1:k,:));
105    pred = mode(pred,1); % takes 1st sorted item for ties
106 end
107
108 %% below code is plotting & data generation - reading it is not necessary in order to
109 %% understand the methods used in this analysis
110
111 %% plotting
112 %% one-time plots
113 % example waveform
114 v=A(:,1)'; Fs = 44100;
115 plot((1:length(v))/Fs,v);
116 xlabel('Time [sec]');
117 ylabel('Amplitude');
118 title('Signal of Interest, v(n)');
119 xlim([0 5]);
120 saveas(gcf,'wavform.png')
121
122 close all;
123 % example spec
124 spectrogram(v);
125 title('Spectrogram of Signal');
126 saveas(gcf,'spec.png')
127
128 %% plots for all cases
129 % scatter, training data
130 close all;
131 scatter(Xtr(1,1:71),Xtr(2,1:71),'^r')
132 hold on
133 scatter(Xtr(1,72:143),Xtr(2,72:143),'og')
134 scatter(Xtr(1,144:216),Xtr(2,144:216),'*b')
135 legend('class 1','class 2','class 3')
136 title('LDA Projection of Case I, Training Data');
137 saveas(gcf,'scat-train-1.png')
138
139 close all;
140 % scatter, testing data
141 scatter(Xts(1,1:18),Xts(2,1:18),'^r')
142 hold on

```



```

143 scatter(Xts(1,19:36),Xts(2,19:36),'og')
144 scatter(Xts(1,37:52),Xts(2,37:52),'*b')
145 legend('class 1','class 2','class 3')
146 title('LDA Projection of Case I, Testing Data');
147 saveas(gcf,'scat-test-1.png')
148
149 close all;
150 % scatter, all data
151 scatter(Xtr(1,1:71),Xtr(2,1:71),'^r')
152 hold on
153 scatter(Xtr(1,72:143),Xtr(2,72:143),'og')
154 scatter(Xtr(1,144:216),Xtr(2,144:216),'*b')
155 scatter(Xts(1,:),Xts(2,),'+k')
156 legend('class 1','class 2','class 3','testing data')
157 title('LDA Projection of Case I, All Data');
158 saveas(gcf,'scat-all-1.png')
159
160 close all;
161 % scatter, errors + training data
162 err_pts = [plab(err.ind); Lts(err.ind)];
163 scatter(Xtr(1,1:71),Xtr(2,1:71),'^r')
164 hold on
165 scatter(Xtr(1,72:143),Xtr(2,72:143),'og')
166 scatter(Xtr(1,144:216),Xtr(2,144:216),'*b')
167 scatter(Xts(1,err.ind),Xts(2,err.ind),'+k');
168 legend('class 1','class 2','class 3','wrongly labelled data')
169 title('LDA Projection of Case I, Errors');
170 saveas(gcf,'scat-err-1.png')
171
172 close all;
173 % energy vs mode
174 lambda = diag(S).^2;
175 semilogy(lambda/sum(lambda),'ko')
176 xlabel('Mode')
177 ylabel('Share of Total Energy')
178 title('Singular Values of Case 1')
179 saveas(gcf,'modes-1.png')
180
181 %% data generation
182 % urls.mat, the script for creating the matrices of urls, is not included here for brevity
183 % (it's over 200 lines long), but it is uploaded to github
184
185 %% get_data.m
186 load('urls.mat')
187
188 % make labels
189 m = zeros(1,90);
190 lab = [m+1 m+2 m+3];
191
192 mFs = getMinFs(Au);
193 A = [];
194 for j=1:size(Au,1)
195     [y,Fs] = webread(Au(j,1));
196     % stereo -> mono
197     if size(y,2)==2
198         y = (y(:,1)+y(:,2))/2;
199     end
200     y = downsample(y, round(Fs/mFs));
201     % take 3 samples, (roughly) 5 sec each
202     A = [A;
203         y(30*mFs:35*mFs,:)';
204         y(45*mFs:50*mFs,:)';
205         y(60*mFs:65*mFs,:)'];
206 end
207
208 mFs = getMinFs(Bu);
209 B = [];
210 for j=87:size(Bu,1)
211     [y,Fs] = webread(Bu(j,1));
212     % stereo -> mono
213     if size(y,2)==2

```

```

214     y = (y(:,1)+y(:,2))/2;
215 end
216 % downsample
217 y = downsample(y, round(Fs/mFs));
218 Fs = Fs/2;
219 % take 3 samples, 5 sec each
220 B = [B;
221     y(30*mFs:35*mFs,:)';
222     y(45*mFs:50*mFs,:)';
223     y(60*mFs:65*mFs,:)'];
224 end
225
226 mFs = getMinFs(Cu);
227 C = [];
228 for j=1:size(Cu,1)
229     [y,Fs] = webread(Cu(j,1));
230     % stereo -> mono
231     if size(y,2)==2
232         y = (y(:,1)+y(:,2))/2;
233     end
234     % downsample
235     y = downsample(y, round(Fs/mFs));
236     Fs = Fs/2;
237     % take 3 samples, 5 sec each
238     C = [C;
239         y(30*mFs:35*mFs,:)';
240         y(45*mFs:50*mFs,:)';
241         y(60*mFs:65*mFs,:)'];
242 end
243
244 A = A'; B = B'; C = C'; % each col is a song now
245 Fs=mFs;
246 save('data.mat','A','B','C','lab','Fs')
247
248 function [minFs] = getMinFs(U)
249     [~,minFs] = webread(U(1));
250     for j=2:size(U,1)
251         [~,Fs] = webread(U(j));
252         if Fs < minFs
253             minFs = Fs;
254         end
255     end
256 end

```