

# AMATH 482 Homework 4

Jonah Wu

March 10, 2021

## Abstract

We are going to use the MNIST data set, which is a collection of handwritten numerical digits (0 - 9), to perform machine learning training to classify different digits using Linear Discriminant Analyzer (LDA), Support Vector Machine (SVM) and Binary Decision Tree with MATLAB. In addition, we are going to find which two digits in the data set appear to be the most difficult to separate, and another pair that are most easy to separate.

## 1 Introduction and Overview

Machine Learning (ML) is one of the hottest topics in the world where there exist many applications of ML to complete tasks that could be accomplished by a human or not by human due to large size of data. In this assignment, Linear Discriminant Analysis (LDA) was used to accomplish the task of training a model that can classify different hand-written digits. Given pictures and labels, Machine Learning is unbiased where it only follows the pattern from the input data set. The two main types of machine learning models are supervised and unsupervised, and in this unbiased case, supervised learning was used to train on the pre-labelled data. In addition, we are going to use and compare 3 different classification models for training the model for classifying digits: linear discriminant analysis, support vector machines, and binary decision trees.

## 2 Theoretical Background

### 2.1 Singular Value Decomposition (SVD)

The SVD is a form of factoring a real or complex matrix that generalizes the eigen decomposition of a square normal matrix to any  $m \times n$  matrix. Singular value decomposition of the matrix  $\mathbf{A}$ :

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

where  $\mathbf{U} \in \mathbb{R}^{m \times m}$  and  $\mathbf{V} \in \mathbb{R}^{n \times n}$  are unitary matrices, and  $\Sigma \in \mathbb{R}^{m \times n}$  is diagonal.

The values  $\sigma_n$  on the diagonal of  $\Sigma$  are called the singular values of the matrix  $\mathbf{A}$ . The number of non-zero singular values is equal to the rank of  $\mathbf{A}$ . The vectors  $u_n$  which make up the columns of  $\mathbf{U}$  are called the left singular vectors of  $\mathbf{A}$ . The vectors  $v_n$  which make up the columns of  $\mathbf{V}$  are called the right singular vectors of  $\mathbf{A}$ .

### 2.2 Linear Discriminant Analysis (LDA)

In order to classify different categories, we use LDA to find a suitable projection that maximizes the distance between the inter-class data while minimizing the intra-class data. We can first find the **between-class scatter matrix**:

$$S_B = (\mu_2 - \mu_1)(\mu_2 - \mu_1)^T \quad (2)$$

where  $\mu$  is the mean of each group. **Between-class scatter matrix** is a measure of the variance between the groups (between the means). Next, we can define the **within-class scatter matrix**:

$$S_w = \sum_{j=1}^2 \sum_x (x - u_j)(x - u_j)^T \quad (3)$$

**Within-class scatter matrix** is a measure of the variance within each group. Finally, our end goal is to find the weight vector  $\mathbf{w}$ :

$$w = \operatorname{argmax} \frac{w^T S_B w}{w^T S_w w} \quad (4)$$

In order to distinguish different labels better, we want to maximize the distance between classes while minimizing distance within class. It turns out that the vector  $\mathbf{w}$  that maximizes the above quotient is the eigenvector corresponding to the largest eigenvalue of the generalized eigenvalue problem:

$$S_B w = \lambda S_w w \quad (5)$$

With this fact, we can easily perform this calculation in MATLAB.

### 3 Algorithm Implementation and Development

With the original file, we will first reshape the data into what we need and then perform the SVD in order to get information about the principal components. After we know which principal components are useful based on the singular value energies, we then use those components to train our model by calculating the between-class scatter matrix and the within-class scatter matrix in order to perform LDA on our data. Finally, with the weighted vectors we trained, we then can test the accuracy on the testing data set. Finally, we will perform SVM and Binary Decision Tree to compare the 3 different classifying methods.

---

#### Algorithm 1: Classifying

---

Import training and testing data
Reshape data
SVD on the training data
Calculate weighted vectors for the model
Use the model to classify the testing data
Compare results with testing labels to find accuracy

---

In the "Calculate weighted vectors for the model" part, LDA is used to compare all pairs among the digits one by one with the between-class scatter matrix and the within-class scatter matrix by using `digits-trainer()` in my MATLAB code. SVM and Binary Decision Tree methods are performed using MATLAB build in functions: `fitcecoc()` and `fitctree()`.

## 4 Computational Results

### 4.1 Singular Value Decomposition

After we perform SVD on our training data set, we will have 3 matrices:  $\mathbf{U}$ ,  $\Sigma$ , and  $\mathbf{V}$ .  $\mathbf{U}$  gives us the left singular vectors in which it represents the principal components.  $\Sigma$  gives us a diagonal with singular values.  $\mathbf{V}$  gives us the right singular values that tells us how each of the images is represented in the Principal Component Analysis basis. Next, we plot the singular value spectrum to help us determine how many modes are necessary for good image reconstruction.

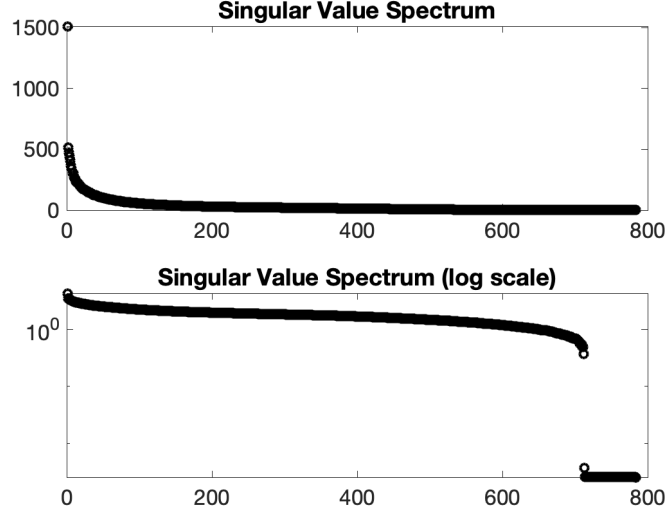


Figure 1: Singular Value Spectrum

As we can see from the top graph, the first principal component is dominating; however, we have almost 800 components to determine so we can scale it in log to know which ones matter. It turns out that after 700th component, there is a significant drop. Although the drop happens at 700, it does not necessarily mean all 700 components are needed. I calculated the energy with each components and found that the first 60 components cover more than 90% of the energy, so I choose 60 as the number of features that are necessary for good image reconstruction in this case.

## 4.2 Projection onto 3 selected V-modes

On a 3D plot, we project onto three selected V-modes (columns 2, 3, and 5) colored by their digit label.

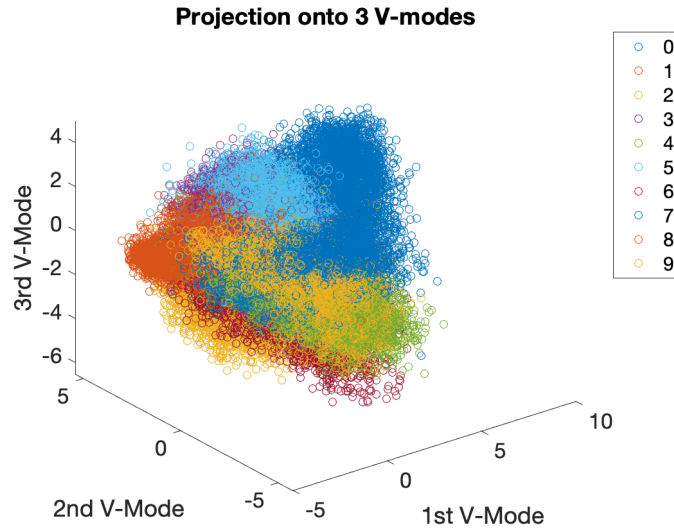
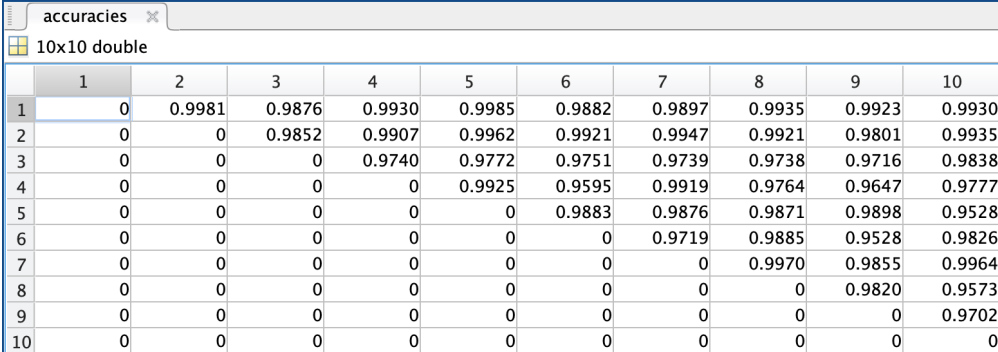


Figure 2: Projection onto 3 V-modes

### 4.3 Linear Discriminant Analysis (LDA)

In my MATLAB code in **Appendix B**, when I am performing LDA on the digits, I have a  $10 \times 10$  matrix **accuracies** that stores the accuracy of each pairs from digits 0 to 9. There is no index 0 for MATLAB so index 1 means 0, and index 10 means 9, etc.



	1	2	3	4	5	6	7	8	9	10
1	0	0.9981	0.9876	0.9930	0.9985	0.9882	0.9897	0.9935	0.9923	0.9930
2	0	0	0.9852	0.9907	0.9962	0.9921	0.9947	0.9921	0.9801	0.9935
3	0	0	0	0.9740	0.9772	0.9751	0.9739	0.9738	0.9716	0.9838
4	0	0	0	0	0.9925	0.9595	0.9919	0.9764	0.9647	0.9777
5	0	0	0	0	0	0.9883	0.9876	0.9871	0.9898	0.9528
6	0	0	0	0	0	0	0.9719	0.9885	0.9528	0.9826
7	0	0	0	0	0	0	0	0.9970	0.9855	0.9964
8	0	0	0	0	0	0	0	0	0.9820	0.9573
9	0	0	0	0	0	0	0	0	0	0.9702
10	0	0	0	0	0	0	0	0	0	0

Figure 3: Accuracy of each pair of digit

As we can see from the above table, row 1 and column 5 has the highest accuracy of 0.9985, and the lowest accuracy of 0.9528 at row 5 and index 10. This tells us that 0 and 4 appear to be the most easy to separate, and 4 and 9 are the most difficult to separate. Overall, LDA has an accuracy of 98.32% on the testing data set.

### 4.4 Support Vector Machines (SVM)

We first calculate the accuracy on the whole data set of 10 digits using SVM and found out that the accuracy is 93.4%. We then calculate its performance on the hardest pair 4 and 9 with an accuracy of 96.38%. while its performance on the easiest pair 0 and 4 is 99.85%

### 4.5 Decision Tree Classifiers

The classification tree has an accuracy of 84.87% when classifying between all 10 digits. For the hardest pair 4 and 9, the tree has an accuracy of 45.16% on the test data for 4 and 9, while having an accuracy of 98.98% on the easiest pair 0 and 4.

## 5 Summary and Conclusions

As we can see the table below, LDA and SVM generally have a better accuracy than the Decision Tree Classifier. In addition, when it comes to the most difficult pair to distinguish: 4 and 9, the Decision Tree Classifier has an accuracy below 50%. Through this assignment, we are able to see applications of machine learning using classifiers such as LDA, SVM, and decision trees. It is really easy for us human to recognize a digit from a hand-written digit image, but what if we need to do this for a large amount of images in a short period of time. We use computers to help us; however, computers do not have eyes and life experience like us to recognize hand-written digits. Machine Learning comes in to help computers to train models with whatever data and labels we give the computers to learn. Overall, we did 3 unbiased supervised machine learning classifier in order to classify different digits with a accuracy of 98.32% (LDA), 93.4% (SVM), and 84.87% (Decision Tree).

accuracy	LDA	SVM	Tree
between all digits	98.32%	93.4%	84.87%
between 0 and 4	99.85%	99.85%	98.98%
between 4 and 9	95.28%	96.38%	45.16%

## Appendix A MATLAB Functions

- `size(A)` returns a row vector whose elements are the lengths of the corresponding dimensions of A.
- `[U, S, V] = svd(X)`: performs a singular value decomposition of matrix A, such that  $A = U \cdot S \cdot V'$ .
- `B = repmat(A,n)` returns an array containing n copies of A in the row and column dimensions. The size of B is  $\text{size}(A) \times n$  when A is a matrix.
- `abs(X)` returns the absolute value of a given input X.
- `diag(X)` returns the diagonal entries of a given a matrix X.
- `eig(X)` returns a column vector of eigenvalues corresponding to the matrix X.
- `fitecoc(X,Y)` returns a trained, multi-class support vector machine, given data X and labels Y.
- `fitctree(X,Y)` returns a trained, multi-class binary decision classification tree, given data X and labels Y.
- `im2double(X)` returns an image converted to double precision, given an image X.
- `mean(X, n)` returns a vector of mean along the n dimension of a given matrix X.
- `norm(X)` returns the Euclidean norm of a vector or matrix X.
- `sort(X)` returns a sorted array of input X, in ascending array.
- `mnist-parse(X)` returns given files downloaded from MNIST.

## Appendix B MATLAB Code

```
1  %% Load data and reshape
2
3  [images, labels] = mnist_parse('train-images-idx3-ubyte', 'train-labels-idx1-ubyte');
4  [t_images, t_labels] = mnist_parse('t10k-images-idx3-ubyte', 't10k-labels-idx1-ubyte');
5
6  training = zeros(784,60000);
7  for i = 1:60000
8      training(:,i) = im2double(reshape(images(:,:,i),784,1));
9  end
10
11 testing = zeros(784,10000);
12 for i = 1:10000
13     testing(:,i) = im2double(reshape(t_images(:,:,i),784,1));
14 end
15
16
17 %% Singular Value Decomposition
18
19 [U,S,V] = svd(training, 'econ');
20
21 figure()
22 subplot(2,1,1)
23 plot(diag(S), 'ko', 'Linewidth',2)
24 title("Singular Value Spectrum")
25 set(gca, 'FontSize',16)
26 subplot(2,1,2)
27 semilogy(diag(S), 'ko', 'Linewidth',2)
28 title("Singular Value Spectrum (log scale)")
29 set(gca, 'FontSize',16)
30 saveas(gcf, 'spectrum.png')
31
```

```

32
33 %% Projection onto 3 V-modes
34
35 projections = U(:, [2,3,5])'*training;
36 for digit=0:9
37     projection = projections(:, labels == digit);
38     plot3(projection(1,:), projection(2,:), projection(3,:), 'o', ...
39         'DisplayName', sprintf('%i', digit))
40     hold on
41 end
42 title('Projection onto 3 V-modes')
43 xlabel('1st V-Mode')
44 ylabel('2nd V-Mode')
45 zlabel('3rd V-Mode')
46 legend
47 set(gca, 'FontSize', 16)
48 saveas(gcf, 'projection.png')
49
50
51 %% Running LDA on all pairs
52
53 feature = 60;
54 accuracies=zeros(10,10);
55 correct = 0;
56 count = 0;
57 for i=0:8
58     for j=i+1:9
59         digit1 = training(:, labels==i);
60         digit2 = training(:, labels==j);
61         [U, w, threshold] = digits_trainer(digit1, digit2, feature);
62
63         test1 = testing(:, t_labels==i);
64         match=0;
65         length1=size(test1,2);
66         for k=1:length1
67             digit = test1(:,k);
68             IMat = U' * digit;
69             digitval = w' * IMat;
70             if digitval < threshold
71                 match = match + 1;
72             end
73         end
74
75         test2 = testing(:, t_labels==j);
76         length2 = size(test2,2);
77         for k=1:size(test2,2)
78             digit = test2(:,k);
79             IMat = U' * digit;
80             digitval = w' * IMat;
81             if digitval > threshold
82                 match = match + 1;
83             end
84         end
85
86         accuracy = match/(length1+length2);
87         accuracies(i+1, j+1) = accuracy;
88         count = count+length1+length2;
89         correct = correct+match;
90     end
91 end
92 LDA_accuracy = correct / count
93
94
95 %% Other Classifiers
96
97 [U, S, V] = svd(training, 'econ');
98 U=U(:, 1:60);
99 projection = S*V';

```

```

100 train = (U'*training)'./max(projection(:));
101 test = (U'*testing)'./max(projection(:));
102
103 % SVM (support vector machines)
104 Mdl = fitcecoc(train,labels);
105 result = predict(Mdl,test);
106 match = result == t.labels;
107 SVM_accuracy = sum(match)/size(match,1)
108
109 % Decision Tree Classifiers
110 d.tree = fitctree(train,labels);
111 result = predict(d.tree,test);
112 match = result == t.labels;
113 tree_accuracy = sum(match)/size(match,1)
114
115
116 %% SVM (support vector machines)
117
118 train_data=train';
119 test_data=test';
120
121 train_0 = train_data(:,labels==0);
122 train_4 = train_data(:,labels==4);
123 train_9 = train_data(:,labels==9);
124
125 test_0 = test_data(:,t.labels==0);
126 test_4 = test_data(:,t.labels==4);
127 test_9 = test_data(:,t.labels==9);
128
129 label_0 = zeros(1,size(train_0,2));
130 label_4 = zeros(1,size(train_4,2)) + 4;
131 label_9 = zeros(1,size(train_9,2)) + 9;
132
133 t_label_0 = zeros(1,size(test_0 ,2));
134 t_label_4 = zeros(1,size(test_4 ,2)) + 4;
135 t_label_9 = zeros(1,size(test_9 ,2)) + 9;
136
137
138 train_04 = [train_0 train_4];
139 test_04 = [test_0 test_4];
140 label_04 = [label_0 label_4];
141 t_label_04 = [t_label_0 t_label_4];
142
143 Mdl_04 = fitcsvm(train_04',label_04);
144 results = predict(Mdl_04,test_04');
145 match = results == t_label_04';
146 SVM_easy_accuracy = sum(match)/size(match,1)
147
148
149 train_49 = [train_4 train_9];
150 test_49 = [test_4 test_9];
151 label_49 = [label_4 label_9];
152 t_label_49 = [t_label_4 t_label_9];
153
154 Mdl_49 = fitcsvm(train_49',label_49);
155 results = predict(Mdl_49,test_49');
156 match = results == t_label_49';
157 SVM_hard_accuracy =sum(match)/size(match,1)
158
159
160 %% Decision Tree Classifiers
161
162 d.tree1 = fitctree(train_04',label_04);
163 results = predict(d.tree1 ,test_04');
164 match = results == t_label_04';
165 tree_easy_accuracy = sum(match)/size(match,1)
166
167 d.tree2 = fitctree(train_49', label_49);

```

```

168 results = predict(d_tree2 ,test_04');
169 match = results == t_label_04';
170 tree_hard_accuracy = sum(match)/size(match,1)
171
172
173 %% digits_trainer function
174
175 function [U,S,V,threshold,w,sort1,sort2] = digits_trainer(d1,d2,feature)
176     n1 = size(d1,2) ;
177     n2 = size(d2,2) ;
178     [U,S,V] = svd([d1 d2], 'econ');
179     digits = S*V';
180     U = U(:,1:feature);
181     digit1 = digits(1:feature,1:n1);
182     digit2 = digits(1:feature,n1+1:n1+n2);
183     ma = mean(digit1,2);
184     mb = mean(digit2,2);
185
186     Sw = 0;
187     for k=1:n1
188         Sw = Sw + (digit1(:,k)-ma)*(digit1(:,k)-ma)';
189     end
190     for k=1:n2
191         Sw = Sw + (digit2(:,k)-mb)*(digit2(:,k)-mb)';
192     end
193     Sb = (ma-mb)*(ma-mb)';
194
195     [V2,D] = eig(Sb,Sw);
196     [~,ind] = max(abs(diag(D)));
197     w = V2(:,ind);
198     w = w / norm(w,2) ;
199     v1 = w'*digit1;
200     v2 = w'*digit2;
201
202     if mean(v1) > mean(v2)
203         w = -w;
204         v1 = -v1;
205         v2 = -v2;
206     end
207
208     sort1 = sort(v1);
209     sort2 = sort(v2);
210     t1 = length(sort1);
211     t2 = 1;
212
213     while sort1(t1) > sort2(t2)
214         t1 = t1-1;
215         t2 = t2+1;
216     end
217     threshold = (sort1(t1) + sort2(t2))/2;
218 end

```



## Appendix C mnist-parse.m Code

```
1 function [images, labels] = mnist_parse(path.to.digits, path.to.labels)
2
3 % The function is courtesy of stackoverflow user rayryeng from Sept. 20,
4 % 2016. Link: ...
5     https://stackoverflow.com/questions/39580926/how-do-i-load-in-the-mnist-digits-and-label-data-in-matlab
6 % Open files
7 fid1 = fopen(path.to.digits, 'r');
8
9 % The labels file
10 fid2 = fopen(path.to.labels, 'r');
11
12 % Read in magic numbers for both files
13 A = fread(fid1, 1, 'uint32');
14 magicNumber1 = swapbytes(uint32(A)); % Should be 2051
15 fprintf('Magic Number - Images: %d\n', magicNumber1);
16
17 A = fread(fid2, 1, 'uint32');
18 magicNumber2 = swapbytes(uint32(A)); % Should be 2049
19 fprintf('Magic Number - Labels: %d\n', magicNumber2);
20
21 % Read in total number of images
22 % Ensure that this number matches with the labels file
23 A = fread(fid1, 1, 'uint32');
24 totalImages = swapbytes(uint32(A));
25 A = fread(fid2, 1, 'uint32');
26 if totalImages ~= swapbytes(uint32(A))
27     error('Total number of images read from images and labels files are not the same');
28 end
29 fprintf('Total number of images: %d\n', totalImages);
30
31 % Read in number of rows
32 A = fread(fid1, 1, 'uint32');
33 numRows = swapbytes(uint32(A));
34
35 % Read in number of columns
36 A = fread(fid1, 1, 'uint32');
37 numCols = swapbytes(uint32(A));
38
39 fprintf('Dimensions of each digit: %d x %d\n', numRows, numCols);
40
41 % For each image, store into an individual slice
42 images = zeros(numRows, numCols, totalImages, 'uint8');
43 for k = 1 : totalImages
44     % Read in numRows*numCols pixels at a time
45     A = fread(fid1, numRows*numCols, 'uint8');
46
47     % Reshape so that it becomes a matrix
48     % We are actually reading this in column major format
49     % so we need to transpose this at the end
50     images(:, :, k) = reshape(uint8(A), numCols, numRows).';
51 end
52
53 % Read in the labels
54 labels = fread(fid2, totalImages, 'uint8');
55
56 % Close the files
57 fclose(fid1);
58 fclose(fid2);
59
60 end
```