

# Classifying Digits

Paul Abboud

**Abstract**—In this paper we analyze the MNIST data set of images by applying Singular Value Decomposition and projecting into Principal Components. Then we performed Linear Discriminant Analysis to construct a linear classifier that can identify digits. Finally, we compare the constructed LDA to Support Vector Machine and Decision Tree classifiers.

## I. INTRODUCTION

The MNIST data set consists of images of digits. In this paper we will perform SVD on the data set to project into Principal Components. Using this information we can determine the number of modes necessary for sufficient image recognition. We will then build Linear Discriminant Analysis (LDA) classifier and determine the accuracy at each digit the model yields when compared to the true values. We then implement Support Vector Machine (SVM) and Decision Tree classifiers to classify the training data. Finally, we examine the identification performance of all three classifiers when compared to the true values to determine the most accurate method.

## II. THEORETICAL BACKGROUND

To understand the issue at hand, the following will discuss the key ideas behind Singular Value Decomposition, Principal Component Analysis, Linear Discriminant Analysis, and Support Vector Machines and Decision Trees.

### A. Singular Value Decomposition

A Singular Value Decomposition decomposes a matrix into three components such that

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

where  $U \in \mathbb{R}^{m \times m}$  and  $V \in \mathbb{R}^{m \times n}$  are unitary matrices and  $\Sigma \in \mathbb{R}^{m \times n}$  is diagonal.  $U$  is a matrix with orthonormal columns,  $V$  is an orthonormal matrix and  $S$  has non-negative singular values in decreasing order along the diagonal. Their product  $A\hat{x}$  represents an orthogonal scaling by  $V^T\hat{x}$  then an axis-aligned scaling by  $S(V^T\hat{x})$  and finally applying resulting coefficients in the orthonormal basis by  $U(S(V^T\hat{x}))$ . In this application, we will subtract the row-wise mean from the data set before computing the SVD in order to center the data. This will make the average value of a pixel 0 across all the images.

### B. Principle Component Analysis

Given a vector  $X$ , we can compute the covariance of the rows with the following formula

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

This is called the covariance matrix which is a very useful tool in minimizing redundancies in the dataset which is the

purpose of PCA. Applying SVD to PCA with the  $1/(n-1)$  factor we have

$$C_X = AA^T = U\Sigma^2U^T \quad (3)$$

Multiplying by  $U^T$  we have

$$Y = U^TX \quad (4)$$

The covariance of  $Y$  can be described in the same fashion as Eq.(2) and Eq.(3) where

$$C_Y = \frac{1}{n-1}YY^T = U^TAA^TU = \Sigma^2 \quad (5)$$

Thus the elements of  $Y$  aren't correlated since the only non-zero elements of  $\Sigma$  are along the diagonal. Now that we have uncorrelated information in  $Y$ , we have essentially reduced the dimension of the data by eliminating redundancies in the original data set  $X$ .

### C. Linear Discriminant Analysis

The purpose of LDA is to find a projection of the data which maximizes inter-class data distance but also minimizes intra-class data. Essentially, for this application, we are maximizing the separation between different digits while minimizing the separation between data points of the same digit. Which can be found by computing the vector  $w$  where

$$w = \underset{w}{\operatorname{argmax}} \frac{w^TS_Bw}{w^TS_ww} \quad (6)$$

where  $S_B$  is the inter-class scatter matrix (measure of variance between classes) and  $S_w$  is the intra-class scatter matrix (measure of variance within classes). The following two sums describe  $S_B$  and  $S_w$  in this application for digits 0 to 9 (10 classifications)

$$S_B = \sum_{j=0}^9 (\mu_j - \mu)(\mu_j - \mu)^T \quad (7)$$

$$S_w = \sum_{j=0}^9 \sum_x (x - \mu_j)(x - \mu_j)^T \quad (8)$$

where  $\mu$  is the overall mean and  $\mu_j$  is the mean of an individual class. Finally, we can compute  $w$  by solving the eigenvalue problem

$$S_Bw = \lambda S_w w \quad (9)$$

Now we have computed a basis for the projection which minimizes intra-class distance and maximizes inter-class distance which will help us classify images of digits.

#### D. Support Vector Machines and Decision Trees

An SVM maps training data to maximize the distance between classifications and assigns new data to the training data group which it is nearest to. It does this by constructing multiple hyperplanes which separate groups of data. The goal is to maximize the distance of the nearest points to the hyperplane. A decision tree evaluates the data point at each node of the tree and determines which branch to follow. Since the classification is discrete in this application (0 to 9), the branches lead to a class label which is a digit in this case.

#### III. ALGORITHM IMPLEMENTATION

To begin, we load the test and training data from the MNIST data set using a MNIST parser and the following

```
1 [train_images, train_labels] =
    mnist_parse('train-images.idx3-ubyte', 'train-labels.idx1-ubyte');
2 [test_images, test_labels] =
    mnist_parse('t10k-images.idx3-ubyte', 't10k-labels.idx1-ubyte');
```

After loading the data and reshaping it, we can then take the row-wise mean and subtract it from the data set. We can then compute the SVD for Principal Component Analysis.

```
1 mu = mean(train_data, 2);
2 train_data = train_data - repmat(mu, 1,
    length(train_data));
3
4 [U, S, V] = svd(train_data, 'econ');
```

The projection is visualized by plotting a 3D scatter plot of the digits depending on three different modes.

```
1 figure(3);
2 colormap jet
3 for i=0:1:9
4     index = find(train_labels == i);
5     scatter3(V(index, 2), V(index, 3), V(
        index, 5), 20, train_labels(index),
        'r');
6     hold on
7 end
8 xlabel('Column 2 of V')
9 ylabel('Column 3 of V')
10 zlabel('Column 5 of V')
11 legend({'0', '1', '2', '3', '4', '5', '6', '7',
    '8', '9'});
```

After computing the SVD, we can develop an LDA classifier to classify the digits from the training data.

```
1 train = proj(:, 2:10);
2 test_t = (U'*test_data)';
3 test = test_t(:, 2:10);
4
5 class = classify(test, train,
    train_labels);
```

```
6
7 err = sum(abs(test_labels - class) > 0);
8 acc_lda = 1 - err / length(test_labels)
```

Using this method we compare pair-wise digits to determine the least and most accurate digit classification pairs of LDA. Then we implement SVM and Decision Tree classifiers using matlab commands to compare their accuracy in digit classification with LDA. The algorithm for building the SVM can be seen in Appendix B. The following builds the Decision Tree and determines the accuracy.

```
1 train = proj(:, 2:10);
2 test_t = (U'*test_data)';
3 test = test_t(:, 2:10);
4 Mdl = fitctree(train, train_labels, '
    OptimizeHyperparameters', 'auto');
5
6 class = predict(Mdl, test);
7
8 err = sum(abs(test_labels - class) > 0);
9 accuracy_ct = 1 - err / length(test_labels)
```

#### IV. COMPUTATIONAL RESULTS

In this application, the SVD has a very intuitive interpretation.  $U$  represents a basis for the coordinates of images since the original data set is made up of image coordinates.  $\Sigma$  helps us determine which components are significant by applying a scalar to  $V^T$  which is an orthonormal matrix of eigenvectors. Thus  $U\Sigma V^T$  is computed by scaling eigenvectors in  $V^T$  by their significance in  $\Sigma$  and these weighted components are transformed by  $U$  into the original coordinate system where the digits are defined. We found that after applying

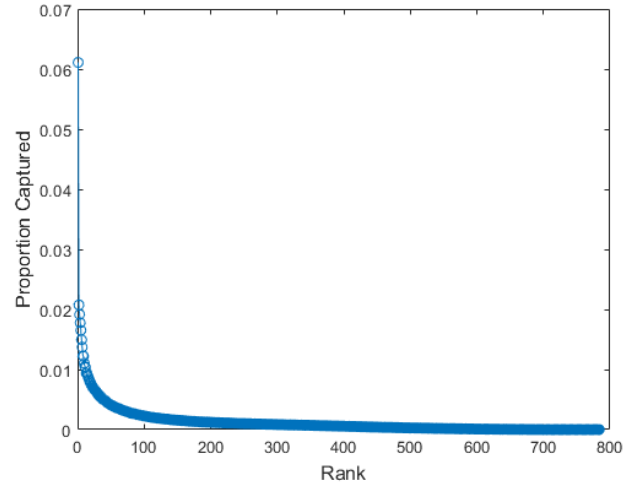


Fig. 1. Energy Captured at Each Dimension in the Image Data

SVD, 90% of the data is captured in 350 modes. This is a significant reduction in dimension when compared to the original data set. We can also see that 60% of the data is

captured in the first rank alone in Figure 1. However, Figure 2 illustrates the reconstruction of digit 5 under different ranks. We can see that it only takes about 100 modes to produce a sufficiently clear image. We can then plot the projection of



Fig. 2. Reconstruction of Digit 5 under various Ranks

three V-modes. Figure 3 shows the projection of the 2nd, 3rd and 5th modes. Each digit is colored and we can see digits clustering which is promising for designing our classifiers on training data. However, there are data points mixing together which could present difficulties and reduce the accuracy of our classifiers. We then implemented LDA and tested the number

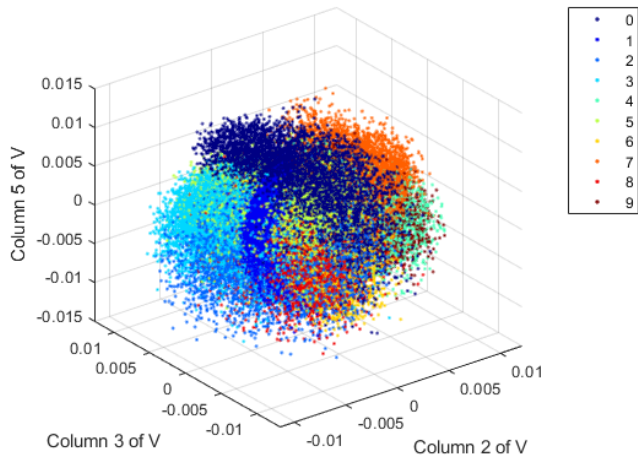


Fig. 3. Projection of modes 2, 3, and 5 and the Resulting Digit Clusters

of V-modes on the classification accuracy. Using digits 0 and 4, we train the LDA classifier and determine the accuracy presented in Figure 4. We can see that at 9 V-modes we have an accuracy of 99.13% which is sufficient. Then we built an LDA for the digits 0, 4 and 7. Figure 5 illustrates the predictions of the constructed LDA. This LDA had an accuracy

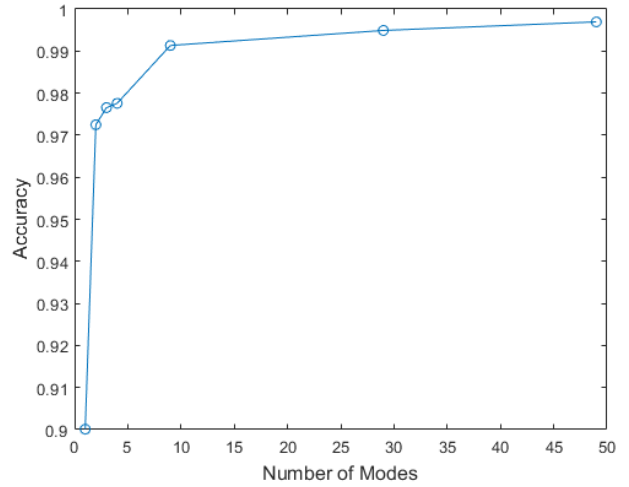


Fig. 4. Two Digit LDA Accuracy

of 96.02% which is highly accurate, but still less accurate when compared to the two digit model. We can identify some inaccuracies from the isolated bars but the majority of digits were classified accurately. Using the trained model, we

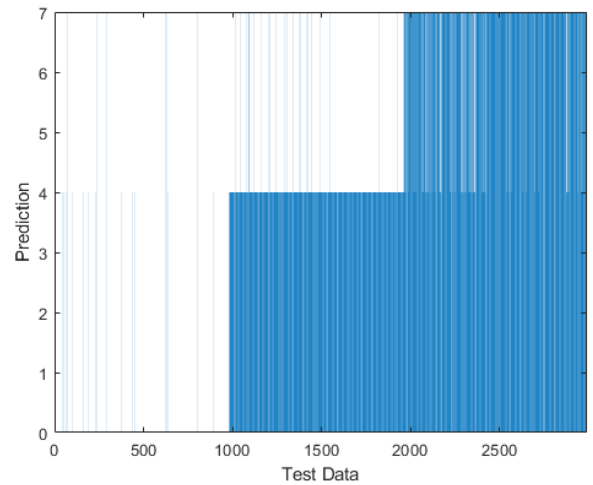


Fig. 5. Three Digit LDA Predictions

determined the accuracy of different digit pairs which is listed in Table 1 and Table 2. Table 1 lists the combinations of digits 0 to 4 with digits 1 to 5. Table 2 lists the combinations of digits 0 to 8 with digits 6 to 9. The digits 4 and 9 had the lowest accuracy of 81.62% meaning 4 and 9 were most difficult to classify. The digits 0 and 1 were easiest to classify with an accuracy of 99.81%. Finally, we built decision tree and SVM classifiers using matlab commands. The decision tree yielded an accuracy of 83.98% while the SVM had an accuracy of 93.32% on the test data. Building an LDA to identify all digits yielded an accuracy of 75.59%, much lower when compared to the two digit classification. Thus the SVM had the highest accuracy when classifying all digits while the

Digit Class	1	2	3	4	5
0	0.9981	0.9732	0.9794	0.9913	0.9498
1		0.9723	0.9828	0.9934	0.9901
2			0.9589	0.9712	0.9574
3				0.9864	0.8985
4					0.9707

TABLE I  
DIGIT PAIR ACCURACY IN PREDICTIONS

Digit Class	6	7	8	9
0	0.9623	0.9895	0.9795	0.9814
1	0.9804	0.9736	0.9445	0.9865
2	0.9266	0.9626	0.9327	0.9691
3	0.9817	0.9637	0.8987	0.9629
4	0.9814	0.9636	0.9693	0.8162
5	0.9492	0.9724	0.8966	0.9548
6		0.9864	0.9731	0.9863
7			0.9401	0.9018
8				0.9319

TABLE II  
DIGIT PAIR ACCURACY IN PREDICTIONS

LDA performed the worst. Finally, we compared the LDA, SVM and decision tree models on the hardest and easiest pair found in the LDA classification. The most difficult pair to classify for LDA was 4 and 9. The least difficult pair to classify was 0 and 1. Table 3 lists the accuracy of each model for the specified pairs. We can see that LDA, SVM and decision tree classifiers were all similarly and significantly accurate in classifying 0 and 1 digits. For the more difficult pair, 4 and 9, we see that SVM performed the best, followed by the decision tree and finally LDA was the least accurate in classifying these digits.

## V. SUMMARY OF RESULTS

In this paper, we applied SVD to significantly reduce the dimension of our data and found that as few as 100 modes could produce identifiable images. We then built an LDA classifier to classify first pairs of digits and then three digits. We found that increasing the number of digit classes lowered the accuracy of the classifier making it inaccurate when classifying all 10 digits with an accuracy of only 75.59%. We then compared the 10 digit LDA classifier to an SVM and decision tree classifier and found that the SVM had the highest accuracy with 93.32%. Finally, we compared the two digit LDA with SVM and decision tree classifiers. We found that all three classifiers had similarly high accuracy when classifying the 0 and 1 digit pair. However, when classifying the 4 and 9 digit pair, we found that the SVM was significantly more accurate than the LDA. Overall, the combination of SVD,

	Classifier Accuracy		
Digit Pairs	LDA	SVM	Decision Tree
0 & 1	0.9981	0.9976	0.9948
4 & 9	0.8162	0.9387	0.8769

TABLE III  
CLASSIFIERS AND THEIR ACCURACY IN DIGIT PAIRS

PCA and LDA significantly reduced the dimension of data and sufficiently classified images of digits, however, the SVM produced the most accurate data.

## APPENDIX A

### MATLAB Functions

**classify:** Classifies the rows of the sample data into groups depending the training data.

**fitcree:** Returns a binary decision tree bases on the input and output matrices. The tree nodes and branches are determined by the data.

**scatter3:** Draws a 3D scatter plot of the data with specified colors.

**fitsvm:** Returns an SVM classifier trained using training data and labels.

**predict:** Returns a vector of predicted class labels for the predicted data.

## APPENDIX B

### MATLAB Code

```

1 close all; clear; clc
2
3 [train_images, train_labels] =
   mnist_parse('train-images.idx3-ubyte', 'train-labels.idx1-ubyte');
4 [test_images, test_labels] =
   mnist_parse('t10k-images.idx3-ubyte', 't10k-labels.idx1-ubyte');
5
6
7 [a,b,c] = size(train_images);
8 [a2,b2,c2] = size(test_images);
9 train_data = zeros(a*b, c);
10 test_data = zeros(a2*b2, c2);
11 for i=1:1:c
12     train_data(:,i)=reshape(
        train_images(:,:,i),a*b, 1);
13 end
14 for i=1:1:c2
15     test_data(:,i)=reshape(test_images
        (:,:,i),a2*b2, 1);
16 end
17
18 mu = mean(train_data,2);
19 train_data = train_data - repmat(mu,1,
        length(train_data));
20
21 [U,S,V]=svd(train_data, 'econ');
22 proj=(S*V)';

```

```

23 plot(diag(S)/sum(diag(S)), '-o')
24 xlabel('Rank','FontSize',12)
25 ylabel('Proportion Captured','FontSize',12)
26
27 figure(3);
28 colormap jet
29 for i=0:1:9
30     index = find(train_labels == i);
31     scatter3(V(index,2),V(index,3),V(index,5),20,train_labels(index),'.')
32     hold on
33 end
34 xlabel('Column 2 of V')
35 ylabel('Column 3 of V')
36 zlabel('Column 5 of V')
37 legend({'0','1','2','3','4','5','6','7','8','9'});
38
39 col_v = [2,3,4,5,10,30,50];
40 acc = zeros(1,length(col_v));
41 for i=1:length(col_v)
42     train_1 = proj(train_labels == 0, 2:col_v(i));
43     train_2 = proj(train_labels == 4, 2:col_v(i));
44     [res_1,y] = size(train_1);
45     [res_2,y] = size(train_2);
46     train = [train_1; train_2];
47     adj_train = [0*ones(res_1,1); 4*ones(res_2,1)];
48
49     test_t = (U'*test_data)';
50     test_1 = test_t(test_labels == 0, 2:col_v(i));
51     test_2 = test_t(test_labels == 4, 2:col_v(i));
52     [res_1,y] = size(test_1);
53     [res_2,y] = size(test_2);
54     test = [test_1; test_2];
55     adj_test = [0*ones(res_1,1); 4*ones(res_2,1)];
56
57     class = classify(test, train, adj_train);
58
59     err = sum(abs(adj_test - class)>0);
60
61     acc(i) = 1-err/length(adj_test);
62 end
63 figure(4)
64 plot(col_v-1, acc, '-o')
65 xlabel('Number of Modes','FontSize',12)

```

```

66 ylabel('Accuracy','FontSize',12)
67
68
69 acc_lda = zeros(10,10);
70 for i=0:1:8
71     for j=i+1:1:9
72         train_1 = proj(train_labels == i,2:10);
73         train_2 = proj(train_labels == j,2:10);
74         [res_1,y] = size(train_1);
75         [res_2,y] = size(train_2);
76         train = [train_1; train_2];
77         adj_train = [i*ones(res_1,1); j*ones(res_2,1)];
78
79         test_t = (U'*test_data)';
80         test_1 = test_t(test_labels==i,2:10);
81         test_2 = test_t(test_labels==j,2:10);
82         [res_1,y] = size(test_1);
83         [res_2,y] = size(test_2);
84         test = [test_1; test_2];
85         adj_test = [i*ones(res_1,1); j*ones(res_2,1)];
86
87         class = classify(test, train, adj_train);
88
89         err = sum(abs(adj_test-class)>0);
90         acc_lda(i+1,j+1)=1-err/length(adj_test);
91     end
92 end
93
94 train_1 = proj(train_labels == 0,2:10);
95 train_2 = proj(train_labels == 4,2:10);
96 train_3 = proj(train_labels == 7,2:10);
97 [res_1,y] = size(train_1);
98 [res_2,y] = size(train_2);
99 [res_3,y] = size(train_3);
100 train=[train_1; train_2; train_3];
101 adj_train = [0*ones(res_1,1); 4*ones(res_2,1); 7*ones(res_3,1)];
102
103 test_t = (U'*test_data)';
104 test_1 = test_t(test_labels == 0,2:10);
105 test_2 = test_t(test_labels == 4,2:10);

```

```

106 test_3 = test_t(test_labels == 7,2:10)
107     ;
108 [res_1,y] = size(test_1);
109 [res_2,y] = size(test_2);
110 [res_3,y] = size(test_3);
111 test = [test_1; test_2; test_3];
112 adj_test = [0*ones(res_1,1); 4*ones(
113     res_2,1); 7*ones(res_3,1)];
114
115 class = classify(test,train,adj_train)
116     ;
117
118 err = sum(abs(adj_test-class)>0);
119 accuracy_047 = 1-err/length(adj_test)
120
121 figure(5)
122 bar(class)
123 xlabel('Test Data')
124 ylabel('Prediction')
125
126 train = proj(:,2:10);
127 test_t = (U'*test_data)';
128 test = test_t(:,2:10);
129
130 class = classify(test,train,
131     train_labels);
132
133 err = sum(abs(test_labels-class)>0);
134 acc_lda = 1-err/length(test_labels)
135
136 train = proj(:,2:10);
137 test_t = (U'*test_data)';
138 test = test_t(:,2:10);
139 Mdl = fitctree(train,train_labels,'
140     OptimizeHyperparameters','auto');
141
142 class = predict(Mdl,test);
143
144 err = sum(abs(test_labels-class)>0);
145 accuracy_ct = 1-err/length(test_labels
146     )
147
148 train = proj(:,2:10)/max(max(S));
149 test = test_t(:,2:10)/max(max(S));
150
151 SVM = cell(10,1);
152 group = 0:1:9;
153 rng(1);
154 for j = 1:numel(group)
155     indx = train_labels==group(j);
156     SVM{j} = fitsvm(train,indx,'
157         ClassNames',[false true],'
158         Standardize',true,...

```

```

154         'KernelFunction','rbf','
155         BoxConstraint',1);
156 end
157 for j = 1:numel(group)
158     [~,score] = predict(SVM{j},test);
159     Scores(:,j) = score(:,2);
160 end
161
162 [~,maxi] = max(Scores,[],2);
163 err = sum(abs(test_labels+1-maxi)>0);
164 accuracy_svm = 1-err/length(
165     test_labels)
166
167 digits = [4,9];
168
169 train_1 = proj(train_labels == digits
170     (1),2:10);
171 train_2 = proj(train_labels == digits
172     (2),2:10);
173 [res_1,y] = size(train_1);
174 [res_2,y] = size(train_2);
175 train = [train_1; train_2];
176 adj_train = [digits(1)*ones(res_1,1);
177     digits(2)*ones(res_2,1)];
178
179 test_t = (U'*test_data)';
180 test_1 = test_t(test_labels == digits
181     (1),2:10);
182 test_2 = test_t(test_labels == digits
183     (2),2:10);
184 [res_1,y] = size(test_1);
185 [res_2,y] = size(test_2);
186 test = [test_1; test_2];
187 adj_test = [digits(1)*ones(res_1,1);
188     digits(2)*ones(res_2,1)];
189
190 rng default
191 Md2 = fitsvm(train,adj_train,'
192     Standardize',true,...
193     'KernelFunction','rbf','
194     BoxConstraint',1);
195
196 class = predict(Md2,test);
197
198 err = sum(abs(adj_test-class)>0);
199 acc_svm_2 = 1-err/length(adj_test)
200
201 Md3 = fitctree(train,adj_train,'
202     OptimizeHyperparameters','auto');
203
204 class = predict(Md3,test);
205
206 err = sum(abs(adj_test-class)>0);
207 accuracy_ct_2 = 1-err/length(adj_test)

```