

Elon Brange, 埃隆

0645120

# Linear Discriminant Analysis (LDA, Fishers discriminant analysis)

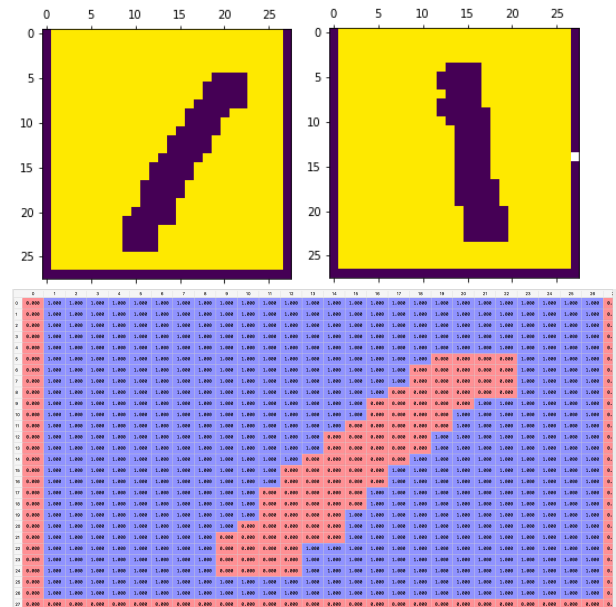
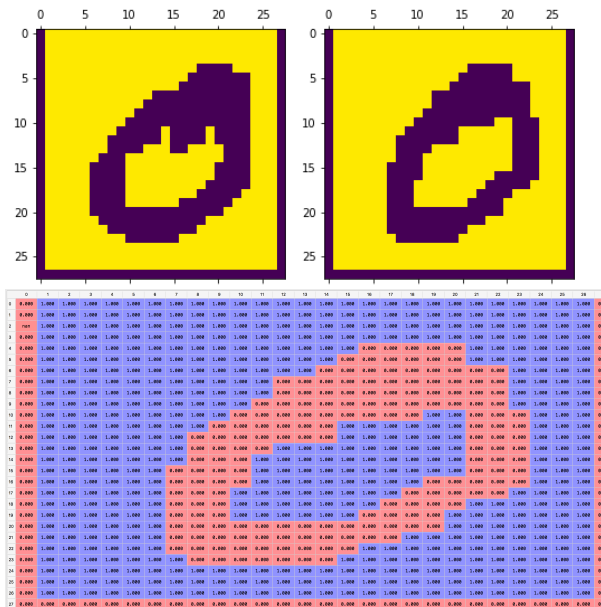
## –Extra Homework, Machine Learning 5255,9i NCTU

### Dataset:

Rows of length 784 representing 28x28 images. Visualization below. (784 features)

**0.txt:** 234 samples

**1.txt:** 277 samples



Visualization code snippet:

```
21 def disp_img(data, n):
22     y = np.empty((28, 28))
23     for i in range(1, 27):
24         low = 28*(i-1)
25         up = 28*i - 1
26         for j in range(0, 27):
27             y[i,j] = data[n, low:up][j]
28     y = np.transpose(y)
29     plt.matshow(y)
30     plt.show
31     return y
```

## LDA – Linear Discriminant Analysis

### 1. Importing dataset

```
16 #Importing dataset
17 data0 = np.genfromtxt("0.txt", delimiter = ",")
18 data1 = np.genfromtxt("1.txt", delimiter = ",")
19
```

### 2. Calculate mean vectors

Mean vectors for each class.

```
45 #Mean vectors Class0 and Class1
46 m0 = np.mean(data0, 0)
47 m1 = np.mean(data1, 0)
48
49 mv = []
50 mv.append(m0)
51 mv.append(m1)
```

$$m_j = \frac{1}{n_j} \sum_{i \in C_j} x_i, \quad j = 1, 2$$

Elon Brange, 埃隆  
0645120

### 3. Scatter matrices

#### a. Within-class scatter

```
52
53 #Within-class scatter matrix
54 #Calculating within-scatter for each class separately
55 #Class0
56 Sw0 = np.zeros((784, 784))
57 for i in range(1, len(data0)):
58     row = data0[i,:].reshape(784,1)
59     mvec = m0.reshape(784, 1)
60     Sw0 += (row - mvec).dot((row - mvec).T)
61
62 #Class1
63 Sw1 = np.zeros((784, 784))
64 for i in range(1, len(data1)):
65     row = data1[i,:].reshape(784,1)
66     mvec = m1.reshape(784, 1)
67     Sw1 += (row - mvec).dot((row - mvec).T)
68
69 #Within-scatter matrix
70 Sw = Sw0 + Sw1
71
```

Calculating within-class scatter for each class:

$$\text{within-class scatter: } s_j^2 = \sum_{i \in \mathcal{C}_j} (y_i - m_j)^2, \quad j = 1, 2$$

Within-class scatter matrix for all data is:  $s_1^2 + s_2^2$

$$\Rightarrow S_w = S_{w \text{ class } 0} + S_{w \text{ class } 1}$$

#### b. Between-class scatter

```
84 #Between-class scatter matrix
85 #Calculating overall mean for each feature
86 mOverall = (m0 + m1)/nclass
87 mOverall = mOverall.reshape(784, 1)
88
89 #B/w-scatter matrix class0
90 mv0 = m0.reshape(784, 1)
91 Sb0 = (mv0 - mOverall).dot((mv0 - mOverall).T)
92
93 #B/w-scatter matrix class1
94 mv1 = m1.reshape(784, 1)
95 Sb1 = (mv1 - mOverall).dot((mv1 - mOverall).T)
96
97 #Between-scatter matrix
98 Sb = 234*Sb0 + 277*Sb1
99
```

First calculate overall mean then proceeding with calculating the between-class scatter matrix:

$$S_b = \sum_{i=1}^{\text{classes}} N_i (\mathbf{m}_i - \mathbf{m}_{tot})(\mathbf{m}_i - \mathbf{m}_{tot})$$

Where,

$N_i$  = sample size

$\mathbf{m}_i$  = sample mean

$\mathbf{m}_{tot}$  = overall mean

### 4. Solving generalized eigenvalue problem

We want to maximize between-class scatter and minimize within-class scatter:

$$J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2}$$

Via Rayleigh quotient and assuming that the within-class scatter matrix ( $S_w$ ) is invertible we derive the matrix:

$$\Rightarrow \mathbf{w} = S_w^{-1}(\mathbf{m}_2 - \mathbf{m}_1) \quad \text{Equals to } S_w^{-1}S_b$$

```
105 #Solving eigenvalue problem
106 #Calculating (Moore-Penrose) pseudo inverse since Sw is singular ->
107 SwInv = linalg.pinv(Sw)
108
109 #Eig to solve
110 eigVal, eigVec = linalg.eig(SwInv.dot(Sb))
111
```

In this case the matrix  $S_w$  is singular and thereby not invertible and I therefore resort to calculate the pseudo inverse.

### 5. Selecting linear discriminants for the new subspace

```
107 #Eigenvalues and Eigenvectors tuples/pairs
108 eigPair = [(np.abs(eigVal[i]), eigVec[:,i]) for i in range(len(eigVal))]
109
110 #Sorting Eigenvalues in descending order
111 eigPair = sorted(eigPair, key=lambda k: k[0], reverse = True)
112
```

Arranging eigenvalues as tuples with their corresponding eigenvectors then sorting them in

descending order. Eigenvectors with the lowest eigenvalue contains the least amount of

Elon Brange, 埃隆

0645120

information and should be dropped. Choosing  $k$  eigenvectors to make a 1-D transformation based on the highest eigenvalue makes us retain as much information as possible.

```
113 #Choosing eigenvectors with largest eigenvalues
114 W = np.hstack((eigPair[0][1].reshape(784,1)))
115
```

```
118 #Explanation of variance
119 varExpl = []
120 eigSum = sum(eigVal)
121 for i,j in enumerate(eigPair):
122     varExpl.append('eigenvalue {0:}: {1:.2%}'.format(i+1, (j[0]/eigSum).real))
123
```

idx ▲	Type	Size	Value
0	str	1	eigenvalue 1: 100.00%
1	str	1	eigenvalue 2: 0.00%
2	str	1	eigenvalue 3: 0.00%
3	str	1	eigenvalue 4: 0.00%
4	str	1	eigenvalue 5: 0.00%
5	str	1	eigenvalue 6: 0.00%
6	str	1	eigenvalue 7: 0.00%
7	str	1	eigenvalue 8: 0.00%

Number of linear discriminants should at most be Classes – 1 and since this is a 2-class problem it should be 1 linear discriminant (eigenvalue = 0). Resulting in only one eigenvector with eigenvalue 0 and we can

form a 1-D feature space based on this eigenpair without losing a lot of information. And as can be seen by looking at the variance the first eigenpair is containing most information and a subspace based on this pair will not loose a lot of information.

## 6. Transforming sample data to the new subspace

```
116 #Transforming data samples to new sub-space
117 lda0 = data0.dot(W)
118 lda1 = data1.dot(W)
119
```

## 7. Visualization of result

```
134 #Plotting histogram of classes
135 plt.hist(lda0, color = "blue")
136 plt.hist(lda1, color = "orange")
137 plt.title("Histogram of classes")
138 plt.show()

131 #Plotting the samples against gaussian noise
132 noise0 = np.random.normal(0, 1, 234)
133 noise1 = np.random.normal(0, 1, 277)
134
135 plt.scatter(lda0, noise0, marker = "x", color = "blue")
136 plt.scatter(lda1, noise1, marker = "x", color = "orange")
137 plt.legend([0,1], bbox_to_anchor=(1.05, 1), loc=2, borderaxespad=0.)
138 plt.xlabel('LDA, W transform')
139 plt.ylabel('Gaussian noise')
140 plt.show()
```

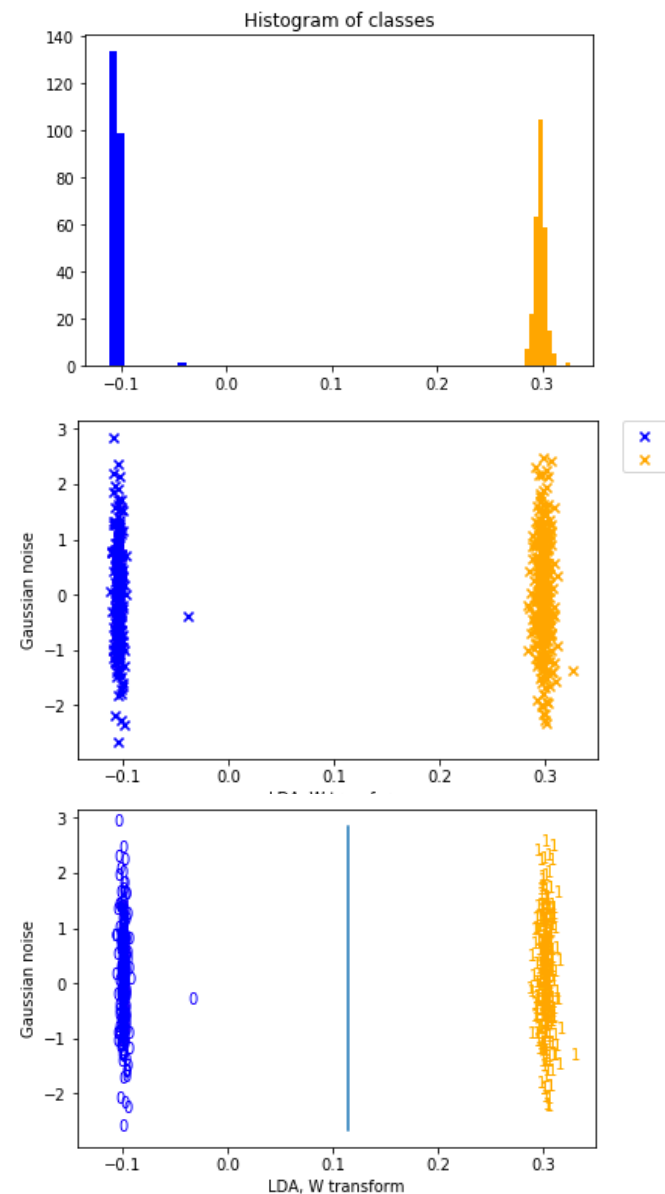
Decision boundary calculated as the hyperplane between the projections of the two means.  
There is no standard rule for calculating the decision boundary.

```
151 #Decision boundary,
152 bl = 0.5*(np.mean(lda1) - np.mean(lda0))
153
154 #Alternative calculation (projecting the means)
155 #bl = 0.5*(m1.dot(W) - m0.dot(W))
156
142 #Alternative plot
143 plt.scatter(lda0, noise0, color = "white")
144 plt.scatter(lda1, noise1, color = "white")
145 for x, y in zip(lda0, noise0):
146     plt.text(x, y, str(0), color="blue", fontsize=10)
147
148 for x, y in zip(lda1, noise1):
149     plt.text(x, y, str(1), color="orange", fontsize=10)
150
151 plt.legend([0,1], bbox_to_anchor=(1.05, 1), loc=2, borderaxespad=0.)
152 plt.xlabel('LDA, W transform')
153 plt.ylabel('Gaussian noise')
154 plt.show()
155
```

## Conclusions and reflections

Implementing the Linear Discriminant Analysis to reduce the dimension of the provided datasets was a quite straightforward process that could, as seen above be broken down into a few steps. These steps are taken from lecture notes provided by professor Wei-Chen Chiu (Slide 180 – 184).

How to choose a proper decision boundary for this dataset seems to be hard. I tried several different approaches but finding a good hyper plane from the equations given above seemed impossible. With the chosen hyper plane lying between the projected means I was not satisfied and decided to see how kmeans performs on the projected dataset since it seems an easy approach and the number of clusters is already known. Further, since I choose the decision boundary to be found as a hyperplane and nothing more complex for this 2-cluster problem regular kmeans seemed to be alright. Result can be seen below.

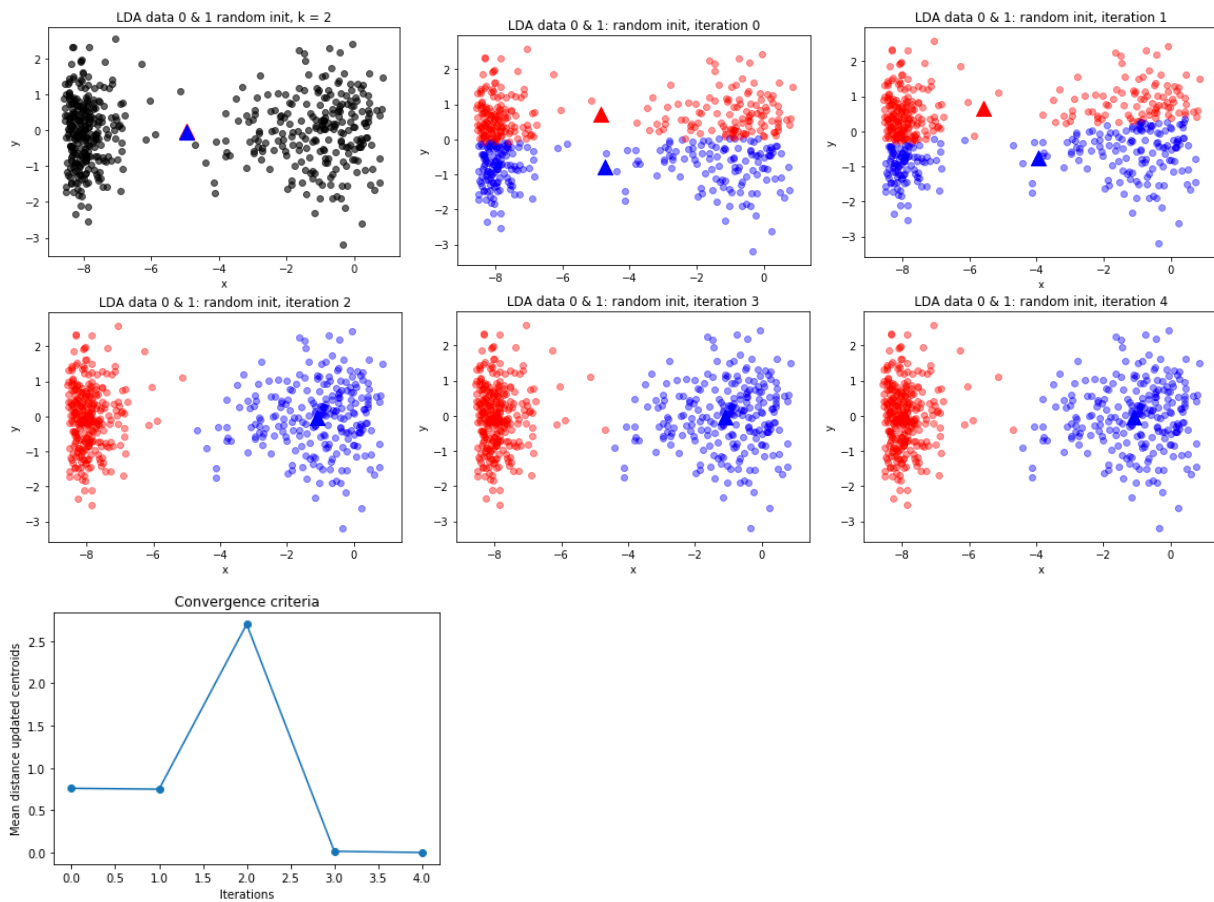


Elon Brange, 埃隆

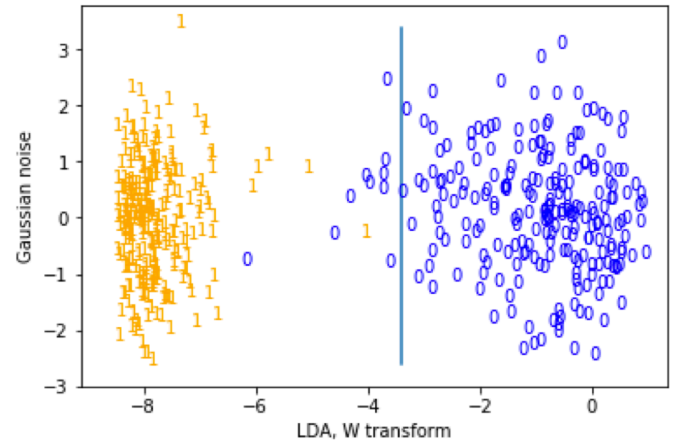
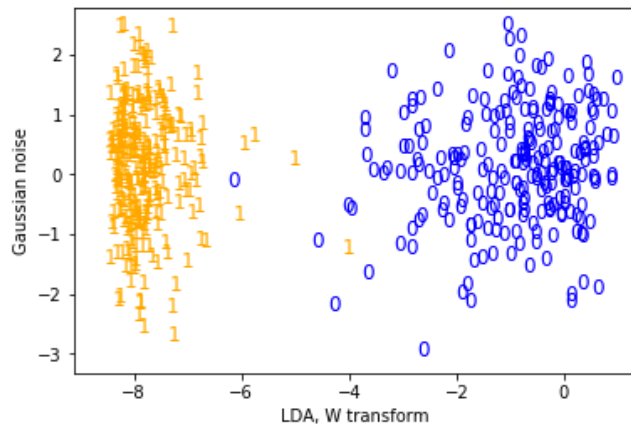
0645120

## K-means clustering of the projected datasets [Not updated...]

Importing the datasets and applying kmeans clustering to classify the dataset results in a correct classification rate of 0.9941291585127201%.



Elon Brange, 埃隆  
0645120



Old inaccurate plots. From wrongly assuming that the matrix was symmetric.