

ML HW2

Part 1.

1. Mean vectors

```
mean vector of class 1: [ 0.99253136 -0.99115481] mean vector of class 2: [-0.9888012  1.00522778]
```

2. Within-class scatter matrix

```
Within-class scatter matrix SW: [[ 2.31560974 -0.95793362]
 [-0.95793362  1.50940428]]
```

3. Between-class scatter matrix

```
Between-class scatter matrix SB: [[ 3.92567873 -3.95549783]
 [-3.95549783  3.98554344]]
```

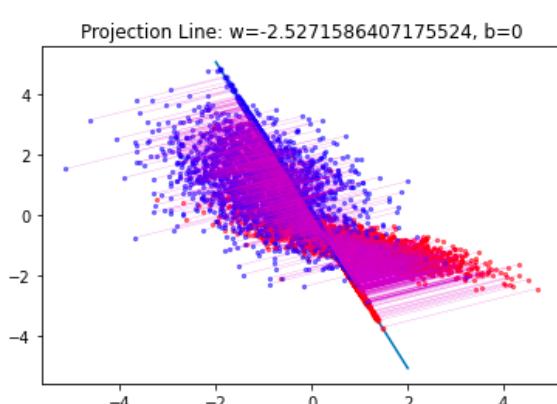
4. Fisher's linear discriminant

```
Fisher's linear discriminant: [-0.36794232  0.92984862]
```

5. Accuracy score on test data with K values from 1 to 5

```
k = 1:
Accuracy of test-set 0.8584
k = 2:
Accuracy of test-set 0.8656
k = 3:
Accuracy of test-set 0.8816
k = 4:
Accuracy of test-set 0.888
k = 5:
Accuracy of test-set 0.8928
```

6. Plotting



Part 2.

What's the difference between the Principle Component Analysis and Fisher's Linear Discriminant?

Principle Component Analysis focuses on capturing the direction of maximum variation in the data set. Fisher's Linear Discriminant focuses on finding a feature subspace that maximizes the separability between the groups.

Please explain in detail how to extend the 2-class FLD into multi-class FLD (the number of classes is greater than two).

When there are only two classes, we add up the two covariance matrices, to extend it to a K-class FLD, the within-class covariance matrix is the sum of all K covariance matrices. As for the between-class covariance matrix, instead of subtracting two means, we get the sum of differences of the mean of each class and a mean in general, also, the rank of this matrix would be at most K-1. The optimal w is the eigenvector of the inverse of the within-class matrix times the between-class matrix, which corresponds to the largest eigenvalue.

By making use of Eq (1) ~ Eq (5), show that the Fisher criterion Eq (6) can be written in the form Eq (7).

$$\begin{aligned}
 J(w) &= \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2} \\
 &= \frac{(w^T(m_2 - m_1))^2}{\sum_{n \in C_1} (y_n - m_1)^2 + \sum_{n \in C_2} (y_n - m_2)^2} \\
 &= \frac{(w^T(m_2 - m_1))(w^T(m_2 - m_1))^T}{\sum_{n \in C_1} (w^T x_n - m_1)^2 + \sum_{n \in C_2} (w^T x_n - m_2)^2} \\
 &= \frac{w^T(m_2 - m_1)(m_2 - m_1)^T w}{\sum_{n \in C_1} (w^T(x_n - m_1))^2 + \sum_{n \in C_2} (w^T(x_n - m_2))^2} \rightarrow S_B = (m_2 - m_1)(m_2 - m_1)^T \\
 &= \frac{w^T(m_2 - m_1)(m_2 - m_1)^T w}{\sum_{n \in C_1} (w^T(x_n - m_1))^2 + \sum_{n \in C_2} (w^T(x_n - m_2))^2} \rightarrow S_W = \sum_{n \in C_1} (x_n - m_1)(x_n - m_1)^T \\
 &= \frac{w^T S_B w}{w^T S_W w + w^T S_W w} = \frac{w^T S_B w}{w^T S_W w} \#
 \end{aligned}$$

Show the derivative of the error function Eq (8) with respect to the activation a_k for an output unit having a logistic sigmoid activation function satisfies Eq (9).

$$\begin{aligned}
 y_k &= \sigma(a_k), \quad \frac{dy}{da_k} = \sigma'(a_k) = \sigma(1-\sigma) \\
 \frac{dE}{da_k} &= -t_k \frac{1}{y_k} (y_k(1-y_k)) + (1-t_k) \frac{1}{1-y_k} (y_k(1-y_k)) \\
 &= (y_k(1-y_k)) \left(-\frac{t_k}{y_k} + \frac{1-t_k}{1-y_k} \right) \\
 &= y_k(1-t_k) - t_k(1-y_k) \\
 &= y_k - t_k \#
 \end{aligned}$$

Show that maximizing likelihood for a multi-class neural network model in which the network outputs have the interpretation $y_k(x, w) = p(t_k=1 | x)$ is equivalent to the minimization of the cross-entropy error function Eq (10).

$$\begin{aligned}
 P(t|X, w) &= \prod_{k=1}^K y_k(x, w)^{t_k} (1-y_k(x, w))^{1-t_k} \\
 E(w) &= -\ln \prod_{n=1}^N P(t|X_n, w) \\
 &= -\ln \prod_{n=1}^N \prod_{k=1}^K y_k(x_n, w)^{t_{nk}} (1-y_k(x_n, w))^{1-t_{nk}} \\
 &= -\sum_{n=1}^N \sum_{k=1}^K \ln (y_k(x_n, w)^{t_{nk}} (1-y_k(x_n, w))^{1-t_{nk}}) \\
 &= -\sum_{n=1}^N \sum_{k=1}^K (t_{nk} \ln y_k(x_n, w) + (1-t_{nk}) \ln (1-y_k(x_n, w))) \\
 &= -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_k(x_n, w)
 \end{aligned}$$