# Machine Learning HW2
## 109550159 李驊恩

## Coding Part

1. the mean vectors mi (i=1, 2) of each 2 classes on training data

```
mean vector of class 1: [ 0.99253136 -0.99115481] mean vector of class 2: [-0.9888012   1.00522778]
Within-class scatter matrix SW: [[ 4337.38546493 -1795.55656547]
 [-1795.55656547  2834.75834886]]
Between-class scatter matrix SB: [[ 3.92567873 -3.95549783]
 [-3.95549783  3.98554344]]
Fisher's linear discriminant: [-0.37003809  0.92901658]
Accuracy of test-set 0.8912
```

2. the within-class scatter matrix SW on training data

```
mean vector of class 1: [ 0.99253136 -0.99115481] mean vector of class 2: [-0.9888012   1.00522778]
Within-class scatter matrix SW: [[ 4337.38546493 -1795.55656547]
 [-1795.55656547  2834.75834886]]
Between-class scatter matrix SB: [[ 3.92567873 -3.95549783]
 [-3.95549783  3.98554344]]
Fisher's linear discriminant: [-0.37003809  0.92901658]
Accuracy of test-set 0.8912
```

3. the between-class scatter matrix SB on training data

```
mean vector of class 1: [ 0.99253136 -0.99115481] mean vector of class 2: [-0.9888012   1.00522778]
Within-class scatter matrix SW: [[ 4337.38546493 -1795.55656547]
 [-1795.55656547  2834.75834886]]
Between-class scatter matrix SB: [[ 3.92567873 -3.95549783]
 [-3.95549783  3.98554344]]
Fisher's linear discriminant: [-0.37003809  0.92901658]
Accuracy of test-set 0.8912
```
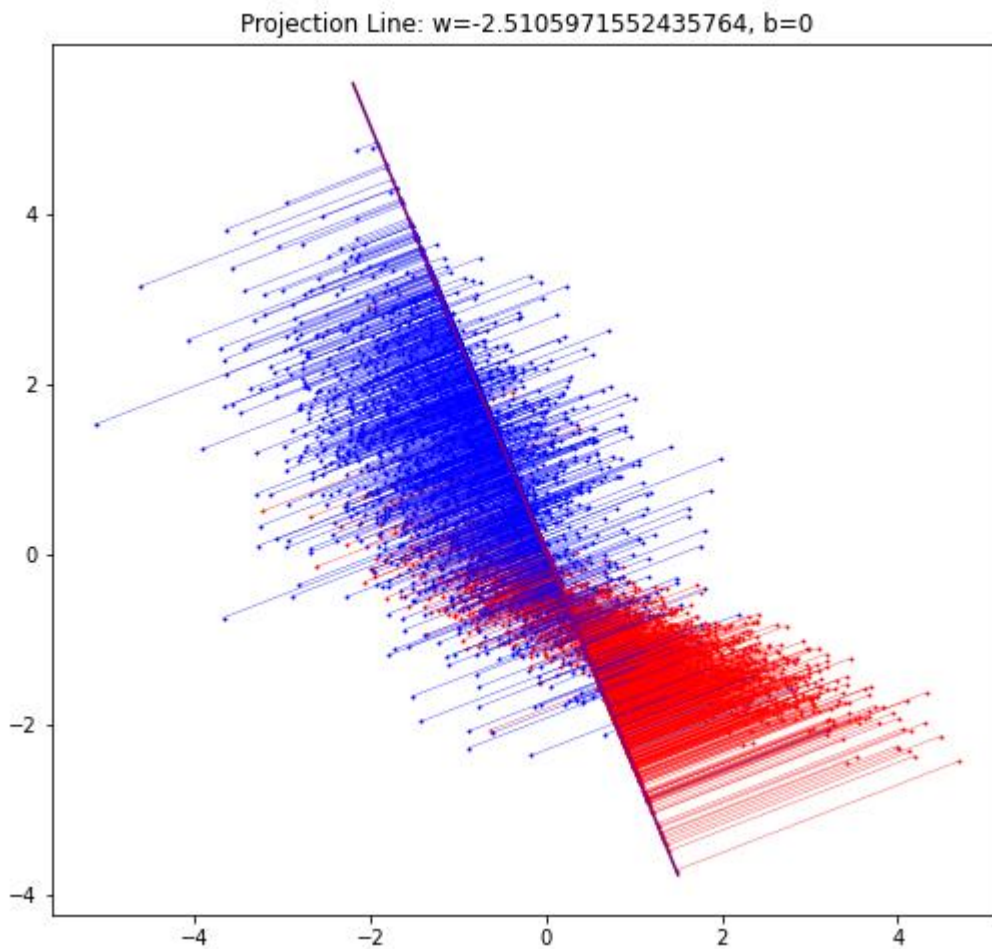
4. the Fisher's linear discriminant w on training data

```
mean vector of class 1: [ 0.99253136 -0.99115481] mean vector of class 2: [-0.9888012   1.00522778]
Within-class scatter matrix SW: [[ 4337.38546493 -1795.55656547]
 [-1795.55656547  2834.75834886]]
Between-class scatter matrix SB: [[ 3.92567873 -3.95549783]
 [-3.95549783  3.98554344]]
Fisher's linear discriminant: [-0.37003809  0.92901658]
Accuracy of test-set 0.8912
```

5. accuracy score on testing data with K values from 1 to 5

K = 5, Accuracy = 0.8912   `Accuracy of test-set 0.8912`

K = 4, Accuracy = 0.8824   `Accuracy of test-set 0.8824`

K = 3, Accuracy = 0.8792   `Accuracy of test-set 0.8792`

K = 2, Accuracy = 0.8704   `Accuracy of test-set 0.8704`

K = 1, Accuracy = 0.8488   `Accuracy of test-set 0.8488`

6.



Projection Line: w=-2.5105971552435764, b=0

# Question Part

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

Principle Component Analysis is an unsupervised technique to reduce dimension. In PCA, we have to find a projection axis to maximize the data separation after projecting. We don't have to know what class the data belong to in PCA. In Fisher's Linear Discriminant, we want the separation of the data to be maximized after projection similarly. What's difference is that we want the separation of the data " between each classes " as large as possible. Thus, FLD is a supervised learning technique since we consider separation between different classes of data, which is different from PCA.

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

2.

To extend the 2-class FLD into multiclass FLD, we first assume that the dimension of input space is $D$, $D > K > 2$

Then linear weight vectors $y_k = W_k^T x$ where $k = 1, \dots D'$, $D' \geq 1$

$\Rightarrow y = W^T x$ where weight vectors $\{W_k\}$ are columns of $W$

The within-class covariance matrix then becomes:

$$S_W = \sum_{k=1}^{K} S_k \quad \text{where } S_k = \sum_{n \in C_k} (X_n - m_k)(X_n - m_k)^T, \quad m_k = \frac{1}{N_k} \sum_{n \in C_k} X_n$$

$N_k$ is the number of pattern in class $C_k$

The between-class covariance matrix then becomes:

$$S_B = \sum_{k=1}^{K} N_k (m_k - m)(m_k - m)^T, \quad m = \frac{1}{N} \sum_{n=1}^{N} X_n$$

However, we can't directly extend the objective that $J(w) = \frac{W^T S_B W}{W^T S_W W}$.

because $W$ and $W^T$ are no longer a single vector. They are matrices that gather the weights vector together and projects each data to a $D'$-dimensional space.

A better choice for objective is $J(w) = \text{Tr}\{(W S_W W^T)^{-1}(W S_B W^T)\}$

As for $D'$, the dimension of the projected space by FLD is at most $K-1$ because the rank of the between-class covariance matrix $S_B$ is at most $K-1$.

3.

By $s_k^2 = \sum_{n \in C_k} (y_n - m_k)^2$, $y = w^T x$, $m_k = w^T m_k$,

$s_1^2 = \sum_{n \in C_1} (y_n - m_1)^2 = \sum_{n \in C_1} (w^T x_n - w^T m_1)^2 = \sum_{n \in C_1} \left[ w^T (x_n - m_1) \right]^T \left[ w^T (x_n - m_1) \right]$

$= \sum_{n \in C_1} \left[ (x_n - m_1)^T w \right]^T \left[ (x_n - m_1)^T w \right] = \sum_{n \in C_1} w^T (x_n - m_1)(x_n - m_1)^T w$

$= w^T \sum_{n \in C_1} (x_n - m_1)(x_n - m_1)^T w$

And thus $s_2^2 = w^T \sum_{n \in C_2} (x_n - m_2)(x_n - m_2)^T w$

Since $S_w = \sum_{n \in C_1} (x_n - m_1)(x_n - m_1)^T + \sum_{n \in C_2} (x_n - m_2)(x_n - m_2)^T$,

$\quad s_1^2 + s_2^2 = w^T S_w w$

By $m_2 - m_1 = w^T (m_2 - m_1)$,

$(m_2 - m_1)^2 = \left[ w^T (m_2 - m_1) \right]^2 = w^T (m_2 - m_1)(m_2 - m_1)^T w$

Since $S_B = (m_2 - m_1)(m_2 - m_1)^T$, $(m_2 - m_1)^2 = w^T S_B w$

By $s_1^2 + s_2^2 = w^T S_w w$ and $(m_2 - m_1)^2 = w^T S_B w$,

$J(w) = \dfrac{(m_2 - m_1)^2}{s_1^2 + s_2^2} = \dfrac{w^T S_B w}{w^T S_w w}$ #

4.

4.

By chain rule, $\dfrac{\partial E}{\partial a_k} = \dfrac{\partial y_k}{\partial a_k} \cdot \dfrac{\partial E}{\partial y_k}$

$$\frac{\partial y_k}{\partial a_k} = \frac{\partial}{\partial a_k}\sigma(a_k) = \frac{\partial}{\partial a_k}\left(\frac{1}{1+e^{-a_k}}\right) = \frac{1}{(1+e^{-a_k})^2} \cdot e^{-a_k} = \frac{1}{1+e^{-a_k}}\left(1 - \frac{1}{1+e^{-a_k}}\right)$$

$$= \sigma(a_k)[1-\sigma(a_k)] = y_k(1-y_k)$$

$$\frac{\partial E}{\partial y_k} = \frac{\partial}{\partial y_k}\Big[-(t_1 \ln y_1 + (1-t_1)\ln(1-y_1)) - (t_2 \ln y_2 + (1-t_2)\ln(1-y_2)) - \cdots$$

$$- (t_k \ln y_k + (1-t_k)\ln(1-y_k)) \cdots - (t_N \ln y_N + (1-t_N)\ln(1-y_N)\Big]$$

$$= 0 + 0 + \cdots + \frac{\partial}{\partial y_k}\Big[-(t_k \ln y_k + (1-t_k)\ln(1-y_k))\Big] + \cdots + 0$$

$$= -\left(\frac{t_k}{y_k} + \frac{1-t_k}{1-y_k}\cdot(-1)\right) = -\frac{t_k}{y_k} + \frac{1-t_k}{1-y_k}$$

$$\Rightarrow \frac{\partial E}{\partial a_k} = \frac{\partial y_k}{\partial a_k}\cdot\frac{\partial E}{\partial y_k} = (y_k(1-y_k))\cdot\left(-\frac{t_k}{y_k} + \frac{1-t_k}{1-y_k}\right) = -(1-y_k)t_k + y_k(1-t_k)$$

$$= -t_k + y_k t_k + y_k - y_k t_k = y_k - t_k$$
$$\#$$

5.

5.

$$y_k(x, w) = p(t_k = 1 | x)$$

$$p(t | w_1, \dots w_k) = \prod_{k=1}^{K} y_k^{t_k}$$

$\Rightarrow$ for data set of N points: $p(T | w_1, \dots w_k) = \prod_{n=1}^{N} \prod_{k=1}^{K} y_{nk}^{t_{nk}}$

To derive the error function, take the negative logarithm:

$$-\ln[p(T | w_1, \dots w_k)] = -\ln \left( \prod_{n=1}^{N} \prod_{k=1}^{K} y_{nk}^{t_{nk}} \right)$$

$$= -\ln \left( \prod_{n=1}^{N} y_{n1}^{t_{n1}} y_{n2}^{t_{n2}} \dots y_{nk}^{t_{nk}} \right)$$

$$= -\ln \left( y_{11}^{t_{11}} \dots y_{1k}^{t_{1k}} \right) \left( y_{21}^{t_{21}} \dots y_{2k}^{t_{2k}} \right) \dots \left( y_{N1}^{t_{N1}} \dots y_{NK}^{t_{NK}} \right)$$

$$= -\left( t_{11} \ln y_{11} + \dots t_{1k} \ln y_{1k} + \dots t_{N1} \ln y_{N1} + \dots t_{NK} \ln y_{NK} \right)$$

$$= -\sum_{n=1}^{N} \sum_{k=1}^{K} t_{nk} \ln y_{nk}$$

$$= -\sum_{n=1}^{N} \sum_{k=1}^{K} t_{nk} \ln y_k(x_n, w) = E(w)$$

\#