

# Intro. to Machine Learning

## Part 1:

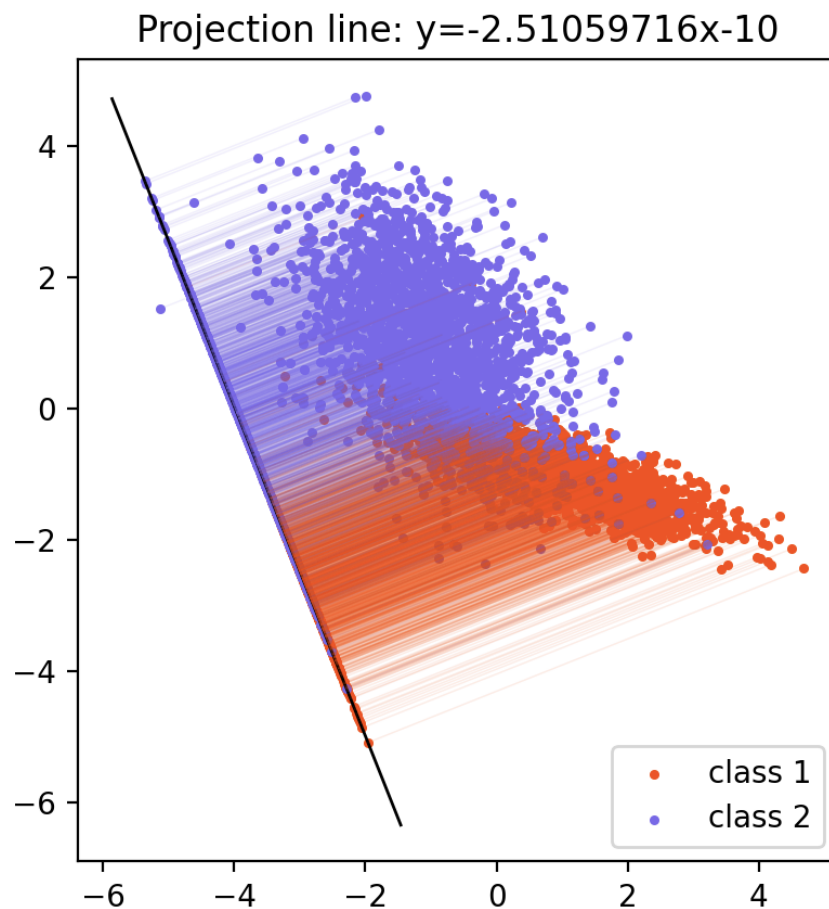
```
mean vector of class 1: [ 0.99253136 -0.99115481] mean vector of class 2: [-0.9888012  1.00522778]

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]]

K=1 accuracy rate: 0.8488
K=2 accuracy rate: 0.8704
K=3 accuracy rate: 0.8792
K=4 accuracy rate: 0.8824
K=5 accuracy rate: 0.8912
```



## Part 2:

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

Principal component analysis is an unsupervised dimensionality reduction method, it ignores the original class label, capturing the direction of maximum variation of the whole dataset; Linear Discriminant is a supervised reduction method that will consider the label while finding the direction maximizing the separability between groups.

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

Assume  $K$ = class number,  $D$ = dimension of input,  $D'$ = dimensional space that weight vector will project to.

The within-class and between-class covariance matrix for multiclass should be:

$$\mathbf{S}_W = \sum_{k=1}^K \mathbf{S}_k \quad \mathbf{S}_B = \sum_{k=1}^K N_k (\mathbf{m}_k - \mathbf{m})(\mathbf{m}_k - \mathbf{m})^T$$

When  $D'=1$ , the objective function:

$$J(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}$$

When  $D'>1$ , the objective function:

$$J(\mathbf{w}) = \text{Tr} \{ (\mathbf{W} \mathbf{S}_W \mathbf{W}^T)^{-1} (\mathbf{W} \mathbf{S}_B \mathbf{W}^T) \}$$

The optimal weight vector is the eigenvector of  $\mathbf{S}_W^{-1} \mathbf{S}_B$  corresponding to the largest eigenvalue, and the dimension of the projected space is at most  $K-1$ .

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

$$y = \mathbf{w}^T \mathbf{x} \quad \text{Eq (1)}$$

$$\mathbf{m}_1 = \frac{1}{N_1} \sum_{n \in \mathcal{C}_1} \mathbf{x}_n \quad \mathbf{m}_2 = \frac{1}{N_2} \sum_{n \in \mathcal{C}_2} \mathbf{x}_n \quad \text{Eq (2)}$$

$$m_2 - m_1 = \mathbf{w}^T (\mathbf{m}_2 - \mathbf{m}_1) \quad \text{Eq (3)}$$

$$m_k = \mathbf{w}^T \mathbf{m}_k \quad \text{Eq (4)}$$

$$s_k^2 = \sum_{n \in \mathcal{C}_k} (y_n - m_k)^2 \quad \text{Eq (5)}$$

$$J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2} \quad \text{Eq (6)}$$

$$J(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} \quad \text{Eq (7)}$$

$$J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2} = \frac{\|\mathbf{w}^T (\mathbf{m}_2 - \mathbf{m}_1)\|^2}{\sum_{n \in \mathcal{C}_1} (\mathbf{w}^T \mathbf{x}_n - m_1)^2 + \sum_{n \in \mathcal{C}_2} (\mathbf{w}^T \mathbf{x}_n - m_2)^2}$$

$$\text{numerator} = [\mathbf{w}^T (\mathbf{m}_2 - \mathbf{m}_1)] [\mathbf{w}^T (\mathbf{m}_2 - \mathbf{m}_1)]^T = \mathbf{w}^T \mathbf{S}_B \mathbf{w}$$

$$\text{denominator} = \sum_{n \in \mathcal{C}_1} (\mathbf{w}^T \mathbf{x}_n - m_1)^2 + \sum_{n \in \mathcal{C}_2} (\mathbf{w}^T \mathbf{x}_n - m_2)^2$$

$$= \mathbf{w}^T \mathbf{S}_{W1} \mathbf{w} + \mathbf{w}^T \mathbf{S}_{W2} \mathbf{w} = \mathbf{w}^T \mathbf{S}_W \mathbf{w}$$

$$J(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}$$

4. 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).

$$E(\mathbf{w}) = - \sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\} \quad \text{Eq (8)}$$

$$\frac{\partial E}{\partial a_k} = y_k - t_k \quad \text{Eq (9)}$$

$$\begin{aligned} y_k &= \sigma(a_k) & \frac{d\sigma}{da} &= \sigma(1-\sigma) \\ \frac{dE(w)}{da_k} &= -t_k \frac{1}{y_k} [y_k(1-y_k)] + (1-t_k) \frac{1}{1-y_k} [y_k(1-y_k)] \\ &= [y_k(1-y_k)] \left[ \frac{1-t_k}{1-y_k} - \frac{t_k}{y_k} \right] \\ &= (1-t_k)y_k - t_k(1-y_k) = y_k - t_k \quad \# \end{aligned}$$

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

$$E(\mathbf{w}) = - \sum_{n=1}^N \sum_{k=1}^K t_{kn} \ln y_k(\mathbf{x}_n, \mathbf{w}) \quad \text{Eq (10)}$$

$$\begin{aligned} p(t|\mathbf{w}) &= \prod_{n=1}^N p(t|\mathbf{x}_n, \mathbf{w}) = \prod_{n=1}^N \prod_{k=1}^K y_k(\mathbf{x}_n, \mathbf{w})^{t_{nk}} \\ \ln p(t|\mathbf{w}) &= \sum_{n=1}^N \sum_{k=1}^K t_{kn} \ln y_k(\mathbf{x}_n, \mathbf{w}) \\ \text{maximizing } \ln p(t|\mathbf{w}) &\text{ is equivalent to} \\ \text{minimizing } -\ln p(t|\mathbf{w}) &= E(\mathbf{w}) \end{aligned}$$