

NYCU Introduction to Machine Learning, Homework 2

109550027 紀竺均

Part. 1, Coding (60%):

1. (5%) Compute the mean vectors m_i ($i=1, 2$) of each 2 classes on training data

```
print(f"mean vector of class 1: {m1}", f"mean vector of class 2: {m2}")
```

✓ 0.6s Python

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

2. (5%) Compute the within-class scatter matrix S_w on training data

```
print(f"Within-class scatter matrix SW: {sw}")
```

✓ 0.4s Python

Within-class scatter matrix SW: [[4337.38546493 -1795.55656547]
[-1795.55656547 2834.75834886]]

3. (5%) Compute the between-class scatter matrix S_B on training data

```
print(f"Between-class scatter matrix SB: {sb}")
```

✓ 0.6s Python

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

4. (5%) Compute the Fisher's linear discriminant W on training data

```
print(f" Fisher's linear discriminant: {w}")
```

✓ 0.4s Python

Fisher's linear discriminant: [-0.000224 0.00056237]

5. (20%) Project the testing data by Fisher's linear discriminant to get the class prediction by K-Nearest-Neighbor rule and report the accuracy score on testing data with K values from 1 to 5 (you should get accuracy over **0.88**)

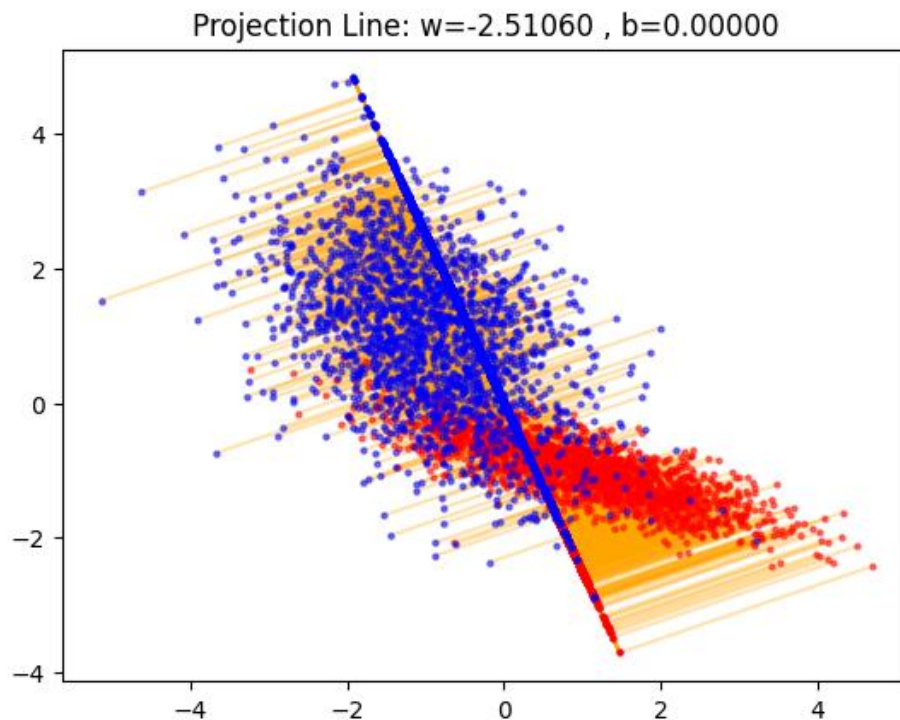
```
print("k = ", k, f"Accuracy of test-set {acc}")
```

✓ 36.4s Python

k = 1 Accuracy of test-set 0.8488
k = 2 Accuracy of test-set 0.8312
k = 3 Accuracy of test-set 0.8792
k = 4 Accuracy of test-set 0.8744
k = 5 Accuracy of test-set 0.8912

6. (20%) Plot the **1) best projection line** on the training data and show the slope and intercept on the title
2) colorize the data with each class

3) project all data points on your projection line.



Part. 2, Questions (40%):

(10%) 1. What's the difference between the Principle Component Analysis and Fisher's Linear Discriminant?

Ans: PCA is an unsupervised dimensionality reduction method while FLD is a supervised technique that takes the class labels into account. Besides, FLD aims to minimize the within class variance and maximize the between class variance.

(10%) 2. Please explain in detail how to extend the 2-class FLD into multi-class FLD (the number of classes is greater than two).

$$\begin{aligned}
 &2. \quad K=2: S_w = S_1 + S_2 \\
 &\quad \rightarrow K>2: S_w = \sum_{k=1}^K S_k, \quad S_k = \sum_{n \in C_k} (x_n - m_k)(x_n - m_k)^T \\
 &\quad K=2: S_B = (m_2 - m_1)(m_2 - m_1)^T \\
 &\quad K>2: S_B = \sum_{k=1}^K N_k (m_k - m)(m_k - m)^T \\
 &\quad (m_k = \frac{1}{N_k} \sum_{n \in C_k} x_n; \quad m = \frac{1}{N} \sum_{n=1}^N x_n) \\
 &\quad K=2: J(w) = \frac{w^T S_B w}{w^T S_w w} \\
 &\quad K>2: J(w) = \frac{w^T S_B w}{w^T S_w w} \\
 &\quad K=2: w = S_w^{-1} (m_2 - m_1) \\
 &\quad K>2: w = \text{eigenvector of } S_w^{-1} S_B \text{ that has Maximal eigenvalue}
 \end{aligned}$$

Ans:

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

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

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

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

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

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

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

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

$$\begin{aligned} 3. J(\mathbf{w}) &= \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2} = \frac{(m_2 - m_1)^2}{\sum_{n \in \mathcal{C}_1} (y_n - m_k)^2 + \sum_{n \in \mathcal{C}_2} (y_n - m_k)^2} \quad (\text{by Eq 5}) \\ &= \frac{(m_2 - m_1)^2}{\sum_{n \in \mathcal{C}_1} (\mathbf{w}^T \mathbf{x}_n - m_k)^2 + \sum_{n \in \mathcal{C}_2} (\mathbf{w}^T \mathbf{x}_n - m_k)^2} \quad (\text{by Eq 1}) \\ &= \frac{[\mathbf{w}^T (m_2 - m_1)]^2}{\sum_{n \in \mathcal{C}_1} (\mathbf{w}^T \mathbf{x}_n - \mathbf{w}^T m_k)^2 + \sum_{n \in \mathcal{C}_2} (\mathbf{w}^T \mathbf{x}_n - \mathbf{w}^T m_k)^2} \quad \begin{matrix} \text{by Eq 3} \\ \rightarrow \text{Eq 4} \end{matrix} \\ &= \frac{(\mathbf{w}^T)^2 (m_2 - m_1)^2}{(\mathbf{w}^T)^2 \cdot \left(\sum_{n \in \mathcal{C}_1} (\mathbf{x}_n - m_1)^2 + \sum_{n \in \mathcal{C}_2} (\mathbf{x}_n - m_2)^2 \right)} \\ &= \frac{(\mathbf{w}^T)^2 \mathbf{S}_B}{(\mathbf{w}^T)^2 \mathbf{S}_W} \quad \left(\frac{1/2 \cdot 1/2 \cdot 1/2 \cdot 1/2}{1/2 \cdot 1/2 \cdot 1/2 \cdot 1/2} \right) \frac{\mathbf{w}}{\mathbf{w}} \\ &= \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}. \end{aligned}$$

Ans:

(7%) 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).

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

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

4. $\frac{\partial E}{\partial a_k} = \frac{\partial E}{\partial y_k} \cdot \frac{\partial y_k}{\partial a_k}$

$y_k = \sigma(a_k) = \frac{1}{1+e^{-a}}$

$\frac{\partial \sigma}{\partial a} = \frac{1}{(1+e^{-a})^2} \cdot (e^{-a}) = \frac{1}{(1+e^{-a})} \cdot \frac{e^{-a}}{(1+e^{-a})} = \sigma(1-\sigma)$

$\Rightarrow \frac{\partial y_k}{\partial a_k} = y_k(1-y_k)$

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

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

$= -t_k(1-y_k) + (1-t_k) \cdot y_k = y_k - t_k$

Ans:

(7%) 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).

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

5. For multiclass neural network:

conditional distribution: $p(t/w_1, w_2, \dots, w_K) = \prod_{k=1}^K y_k^{t_k}$

\rightarrow likelihood function: $p(T/w_1, w_2, \dots, w_K) = \prod_{n=1}^N \prod_{k=1}^K y_{nk}^{t_{nk}}$ (For N points)

\rightarrow negative log likelihood:

$-\ln(P(T/w_1, \dots, w_K)) = -\ln\left(\prod_{n=1}^N \prod_{k=1}^K y_{nk}^{t_{nk}}\right)$

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

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

Ans: