

NYCU Introduction to Machine Learning, Homework 2

109550164 徐聖哲

Part. 2, Questions (40%):.

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

Ans: Both Principle Component Analysis and Fisher's Linear Discriminant are used to find the best linear combination to explain the data. However Fisher's Linear Discriminant is supervised learning, it uses the label of train data to check the new projected dot is classified correct or not. Principle Component Analysis is unsupervised learning, which focus on maximum variance between data.

(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).

Ans: First S_W is the sum of all S_k for $k = k$ class, Then S_B will multiply with number of dot in the class in each between-class variance. Finally, in order to find w , we use lagrangian function to minimize within-class-covariance, thus w is the eigenvector of $(S_W)^{-1} * S_B$ that corresponds to the largest eigenvalue.

2.

①

$$S_W$$
$$k=2: S_W = \sum_{n \in C_1} (x_n - m_1)(x_n - m_1)^T + \sum_{n \in C_2} (x_n - m_2)(x_n - m_2)^T$$

$$k=n: S_W = \sum_{k=1}^K S_k, \quad S_k = \sum_{n \in C_k} (x_n - m_k)(x_n - m_k)^T$$

②

$$S_B$$
$$k=2: S_B = (m_2 - m_1)(m_2 - m_1)^T$$

$$k=n: S_B = \sum_{k=1}^K N_k (m_k - m)(m_k - m)^T$$

③

$$W$$
$$k=2: W = S_W^{-1} (m_2 - m_1)$$

$k=n$: use lagrangian function

$$L_P = -\frac{1}{2} W^T S_B W + \frac{1}{2} \lambda (W^T S_W W - 1)$$

$$\Rightarrow S_B W = \lambda S_W W$$

$$\Rightarrow S_W^{-1} S_B = \lambda W$$

\Rightarrow select optimal w

$w =$ eigenvectors of $S_W^{-1} S_B$ with largest eigenvalue

(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).

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

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

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

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

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

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

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

$$3. J(w) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2}$$

$$\begin{aligned} J_1 &= (m_2 - m_1)^2 = (w^T(m_2 - m_1))^2 \\ &= w^T(m_2 - m_1) (w^T(m_2 - m_1))^T \\ &= w^T(m_2 - m_1) (m_2 - m_1)^T w \\ &= w^T S_b w \end{aligned}$$

$$J_2 = s_1^2 + s_2^2$$

$$\begin{aligned} &= \sum_{n \in C_1} (y_n - m_1)^2 + \sum_{n \in C_2} (y_n - m_2)^2 \\ &= \sum_{n \in C_1} (w^T x_n - w^T m_1)^2 + \sum_{n \in C_2} (w^T x_n - w^T m_2)^2 \\ &= \sum_{n \in C_1} (w^T (x_n - m_1))^2 + \sum_{n \in C_2} (w^T (x_n - m_2))^2 \\ &= w^T \sum_{n \in C_1} (x_n - m_1)^2 w + w^T \sum_{n \in C_2} (x_n - m_2)^2 w \\ &= w^T \left(\sum_{n \in C_1} (x_n - m_1)^2 + \sum_{n \in C_2} (x_n - m_2)^2 \right) w \end{aligned}$$

$$\Rightarrow J(w) = \frac{w^T S_B w}{w^T S_W w}$$

(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).

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

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

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

$$4. \quad y = G(a) = \frac{1}{1 + e^{-a}}$$

$$\star \frac{\partial E}{\partial a_k} = \frac{\partial E}{\partial y} \cdot \frac{\partial y}{\partial a_k}$$

$$\frac{\partial E}{\partial y} = \frac{\partial}{\partial y} (-t_n \ln y_n + (1 - t_n) \ln(1 - y_n))$$

$$= -\frac{t_n}{y_n} + \frac{1 - t_n}{1 - y_n} = \frac{-t_n + t_n y_n + y_n - y_n t_n}{y_n(1 - y_n)} = \frac{y_n - t_n}{y_n(1 - y_n)}$$

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

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

(7%) 5. Show that maximizing likelihood for a multiclass 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).

$$5. y_k(x, w) = p(t_k = 1 | w)$$

$$w = \operatorname{argmax} P(X|w) = \operatorname{argmax} \prod_{i=1}^N \prod_{k=1}^K y_k(x_i, w)^{t_{ik}}$$

then we take negative log-likelihood, turn the problem to minimization and log is strictly increased function

$$\Rightarrow w = \operatorname{argmin} (-\log P(X|w)) = \operatorname{argmin} \left(-\sum_{i=1}^N \sum_{k=1}^K t_{ik} \ln y_k(x_i, w) \right)$$

\Rightarrow maximum likelihood equals minimization of cross entropy error function

$$E(w) = -\sum_{i=1}^N \sum_{k=1}^K t_{ik} \ln y_k(x_i, w)$$