

Q1

From $m_2 - m_1 = w^T(m_2 - m_1)$ we can construct the Lagrangian function

$$L = w^T(m_2 - m_1) + \lambda(w^Tw - 1)$$

Taking the gradient of L we obtain

$$\nabla L = m_2 - m_1 + 2\lambda w.$$

and setting this gradient to zero gives

$$w = -\frac{1}{2\lambda}(m_2 - m_1)$$

from which follows that $w \propto m_2 - m_1$.

Q2. Starting with the numerator on the r.h.s of $J(w)$, we can rewrite it as follows

$$\begin{aligned} (m_2 - m_1)^2 &= (w^T(m_2 - m_1))^2 \\ &= w^T(m_2 - m_1)(m_2 - m_1)^Tw \\ &= w^T S_B w. \end{aligned}$$

Similarly, we can rewrite the denominator of the r.h.s of $J(w)$:

$$\begin{aligned} S_1^2 + S_2^2 &= \sum_{n \in C_1} (y_n - m_1)^2 + \sum_{k \in C_2} (y_k - m_2)^2 \\ &= \sum_{n \in C_1} (w^T(x_n - m_1))^2 + \sum_{k \in C_2} (w^T(x_k - m_2))^2 \\ &= \sum_{n \in C_1} w^T(x_n - m_1)(x_n - m_1)^Tw + \sum_{k \in C_2} w^T(x_k - m_2)(x_k - m_2)^Tw \\ &= w^T S_W w. \end{aligned}$$

Q3. The likelihood function is given by

$$p(\{\phi_n, t_n\} | \{x_k\}) = \prod_{n=1}^N \prod_{k=1}^K \{p(\phi_n | c_k) x_k\}^{t_{nk}}$$

and taking the logarithm, we obtain

$$\ln p(\{\phi_n, t_n\} | \{x_k\}) = \sum_{n=1}^N \sum_{k=1}^K t_{nk} \{\ln p(\phi_n | c_k) + \ln x_k\}$$

In order to maximize the log likelihood with respect to x_k we need to preserve the constraint $\sum_k x_k = 1$. This can be done by introducing a Lagrange multiplier λ and maximizing

$$\ln p(\{\phi_n, t_n\} | \{x_k\}) + \lambda \left(\sum_{k=1}^K x_k - 1 \right).$$

Setting the derivative with respect to x_k equal to zero, we obtain

$$\sum_{n=1}^N \frac{t_{nk}}{x_k} + \lambda = 0.$$

re-arrange the gives.

$$-\sum_k \lambda = \sum_{n=1}^N t_n \lambda = N \lambda.$$

Summing both sides over k we find that $\lambda = -N$, and using this to eliminate λ , we obtain.

$$\lambda_k = \frac{N_k}{N}.$$

Q4. Differentiating $\sigma(a)$ we obtain

$$\frac{d\sigma}{da} = \frac{e^{-a}}{(1-e^{-a})^2} = \sigma(a) \left\{ \frac{e^{-a}}{1+e^{-a}} \right\} = \sigma(a) \left\{ \frac{1+e^{-a}}{1+e^{-a}} - \frac{1}{1+e^{-a}} \right\} = \sigma(a) (1 - \cancel{\sigma(a)})$$

Q5. We start by computing the derivative of $E(w)$ w.r.t y_n .

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

From the result $\frac{ds}{da} = \sigma(1-\sigma)$,

we see that

$$\frac{\partial y_n}{\partial a_n} = \frac{\partial \sigma(a_n)}{\partial a_n} = \sigma(a_n)(1-\sigma(a_n)) = y_n(1-y_n).$$

Finally we have.

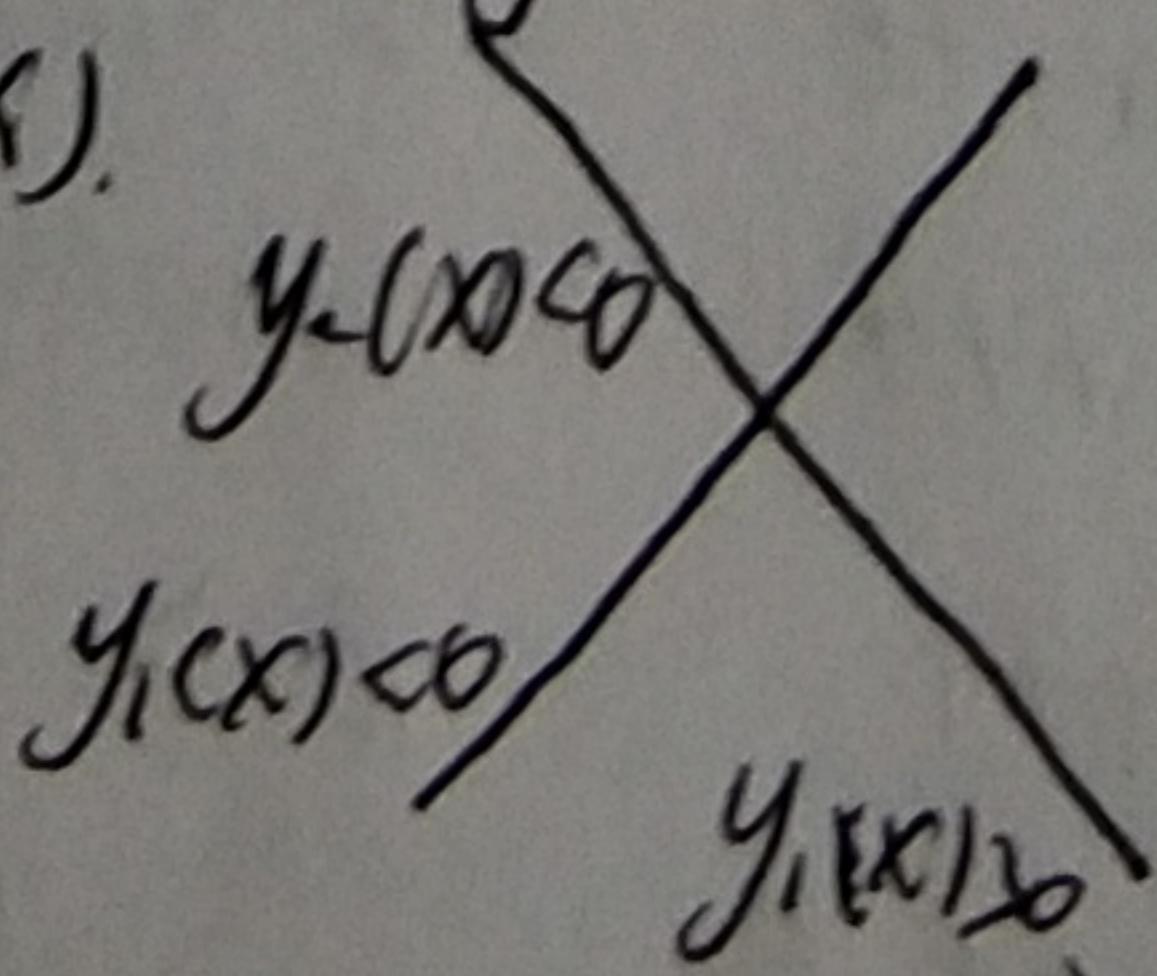
$$\nabla g_n = \phi_n.$$

where ∇ denotes the gradient with respect to w . Combining these three equations using the chain rule, we obtain

$$\nabla E = \sum_{n=1}^N \frac{\partial E}{\partial y_n} \frac{\partial y_n}{\partial a_n} \nabla a_n = \sum_{n=1}^N (y_n - t_n) \phi_n.$$

Q6. P1. Since there are 3 classes, $c=3$ and there are 2 discriminant functions $y_1(x)$ and $y_2(x)$.

For $x \in C_1$, $y_1(x) > 0$, and for $x \in C_2$, $y_2(x) > 0$. This leads to the following problems. How do we classify input patterns x which have the property that $y_1(x) > 0$ and $y_2(x) > 0$? Clearly, they belong to both class C_1 and C_2 . Figure 1 illustrates the problem. Making the discriminant lines parallel to each other does not resolve the problem since the intersection $y_1(x) > 0$ and $y_2(x) > 0$ is non-empty. Note that the intersection $y_1(x) > 0$ and $y_2(x) > 0$ is a null set if and only if the two lines coincide which means $y_1(x) = y_2(x)$.



P2. Since there are 3 classes, $c(c-1)/2 = 3$ and there are three discriminant functions, $y_1(x)$, $y_2(x)$ and $y_3(x)$. The classification structure is as follows.

1. If $y_{12}(x) > 0$ and $y_{13}(x) > 0$, then $x \in C_1$.

2. If $y_{12}(x) < 0$, and $y_{13}(x) > 0$, then $x \in C_2$

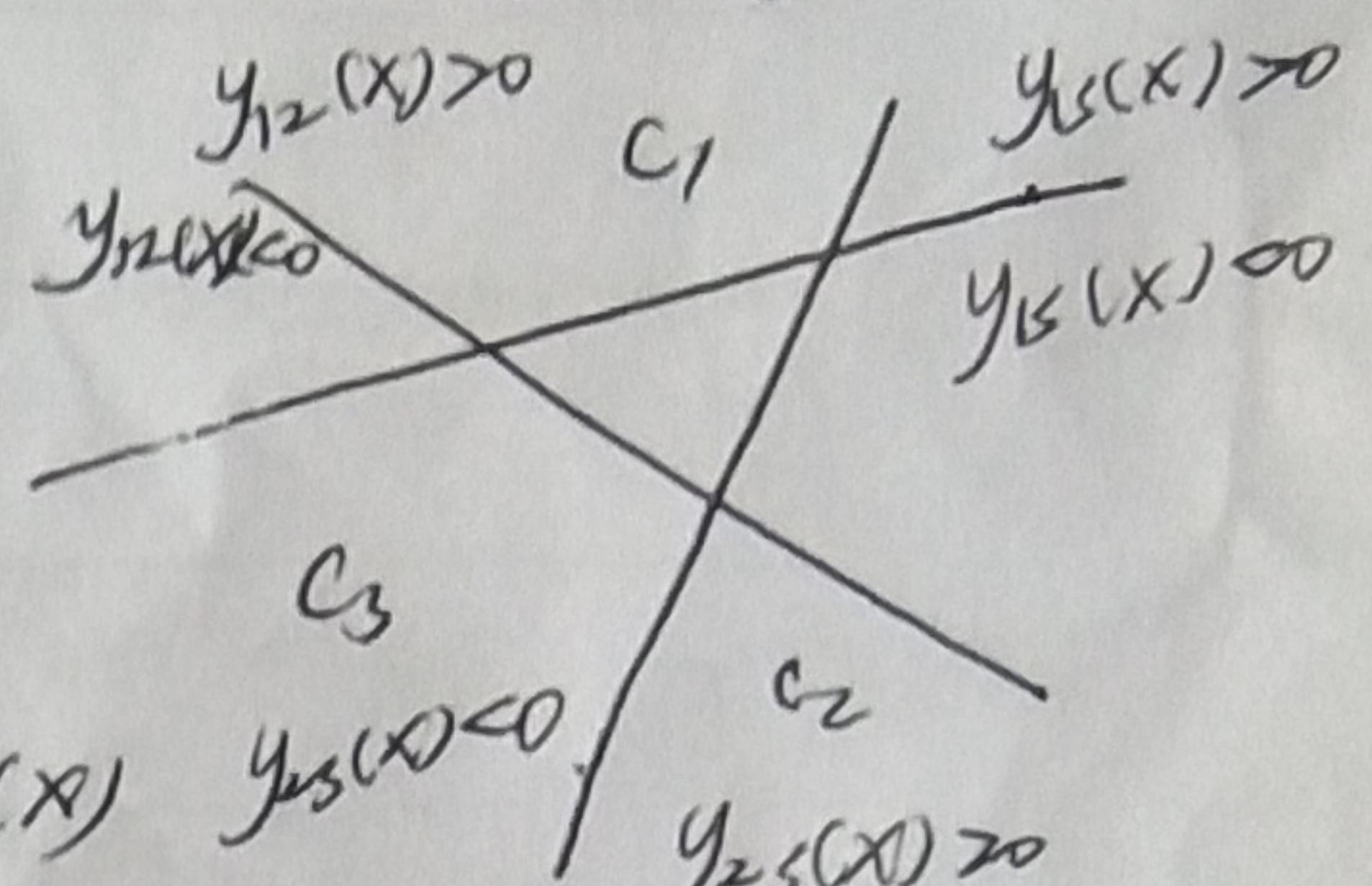
3. If $y_{12}(x) < 0$, and $y_{13}(x) < 0$, then $x \in C_3$

This leads to the following problems as illustrated in figure. The following regions are unclassified.

1. $y_{12}(x) < 0$ and $y_{13}(x) > 0$.

2. $y_{12}(x) > 0$ and $y_{13}(x) > 0$ and $y_{13}(x) < 0$.

The intersections are null sets if and only if $y_{12}(x) = y_{13}(x) = y_{23}(x)$ $y_{12}(x) < 0$, $y_{23}(x) > 0$



Q7. First, let's calculate the linear discriminants for the points belonging to the two convex hulls. For points in the convex hull of $\{x^n\}$, the linear discriminant is:

$$y(x) = \hat{w}^T x^n + w_0. \quad (2)$$

Substituting (1) in (2), we get

$$y(x) = \hat{w}^T (\sum_n \alpha_n x^n) + w_0 \quad (3)$$

Since α_n is a scalar quantity, we can bring the summation in (3) outside \rightarrow

$$y(x) = \sum_n \alpha_n (\hat{w}^T x^n) + w_0 = \sum_n \alpha_n (\hat{w}^T x^n + w_0) \quad (4)$$

where we made use of the fact that $\sum \alpha_n = 1$. Similarly, we can develop the linear discriminant for the ~~problem~~ points belonging to the convex hull of $\{z^m\}$

$$y(z) = \sum_m \beta_m (\hat{w}^T z^m + w_0) \quad (5) \text{ where } \beta_m \geq 0 \text{ and } \sum_m \beta_m = 1$$

Convex hulls intersect: If the convex hulls intersect, there must be at least one point in common, between $\{x^n\}$ and $\{z^m\}$. Let's call that point xz . Since xz belongs to both convex hulls there must be a set of $\{\alpha_n\}$ and $\{\beta_m\}$ that give rise to xz . The linear discriminant for xz can now be written in two separate but equivalent ways. From (4) and (5), we get

$$y(xz) = \sum_n \alpha_n (\hat{w}^T x^n + w_0) = \sum_m \beta_m (\hat{w}^T z^m + w_0) \quad (6)$$

For linear separability, we must have

$$y(x^n) = \hat{w}^T x^n + w_0 > 0 \quad y(z^m) = \hat{w}^T z^m + w_0 < 0. \quad (7)$$

From the non-negativity and simplex constraints on α and β , (6) and (7), we have a contradiction. The linear discriminant $y(xz)$ has to be simultaneously greater than and less than zero, which is impossible.

~~Pattern~~ Patterns are linearly separable:

If the patterns are linearly separable, we know that

$$y(x^n) = \hat{w}^T x^n + w_0 > 0.$$

$$y(z^m) = \hat{w}^T z^m + w_0 < 0 \quad (8)$$

Assume that there is a point x^z lying in the intersection of the convex hulls. From (6) above.

$$y(x^z) = \sum_n \alpha_n (\hat{w}^T x^n + w_0) = \sum_m \beta_m (\hat{w}^T z^m + w_0) \quad (9)$$

The equality in (9) is not possible given the fact from (8) that the patterns are linearly separable.