

# Machine Learning HW5

---

## Q1

We first write down the negative logarithm of the likelihood function

$$E(w, \Sigma) = \frac{1}{2} \sum_{n=1}^N \{[y(x_n, w) - t_n]^T \Sigma^{-1} [y(x_n, w) - t_n]\} + \text{const} \quad (1)$$

The  $\Sigma$  is unknown and  $\text{const}$  denotes the term independent of both  $w$  and  $\Sigma$ . In the first situation, if  $\Sigma$  is fixed and known, the equation above will reduce to:

$$E(w) = \frac{1}{2} \sum_{n=1}^N \{[y(x_n, w) - t_n]^T \Sigma^{-1} [y(x_n, w) - t_n]\} + \text{const} \quad (2)$$

We can simply solve  $w_{ML}$  by minimizing it. If  $\Sigma$  is unknown, since  $\Sigma$  is in the first term on the right of (\*), solving  $w_{ML}$  will involve  $\Sigma$ . Note that in the previous problem, the main reason that they can decouple is due to the independence assumption, i.e.,  $\Sigma$  reduces to  $\beta^{-1}I$ , so that we can bring  $\beta$  to the front and view it as a fixed multiplying factor when solving wML.

## Q2

We know that the logistic sigmoid function  $\sigma(a) \in [0, 1]$ , therefore if we perform a linear transformation  $h(a) = 2\sigma(a) - 1$ , we can find a mapping function  $h(a)$  from  $(-\infty, +\infty)$  to  $[-1, 1]$ . In this case, the conditional distribution of targets given inputs can be similarly written as:

$$p(t|x, w) = \left[ \frac{1 + y(x, w)}{2} \right]^{(1+t)/2} \left[ \frac{1 - y(x, w)}{2} \right]^{(1-t)/2} \quad (3)$$

Where  $[1 + y(x, w)]/2$  represents the conditional probability  $p(C_1|x)$ . Since now  $y(x, w) \in [-1, 1]$ , we also need to perform the linear transformation to make it satisfy the constraint for probability. Then we can further obtain:

$$\begin{aligned}
E(w) &= - \sum_{n=1}^N \left\{ \frac{1+t_n}{2} \ln \frac{1+y_n}{2} + \frac{1-t_n}{2} \ln \frac{1-y_n}{2} \right\} \\
&= -\frac{1}{2} \sum_{n=1}^N \{ (1+t_n) \ln(1+y_n) + (1-t_n) \ln(1-y_n) \} + N \ln 2
\end{aligned} \tag{4}$$

**Q3**

(a)

$$E[t] = \int t N(t|\mu, \sigma^2 I) dt = \mu \tag{5}$$

And

$$E[\|t\|^2] = \int \|t\|^2 N(t|\mu, \sigma^2 I) dt = L\sigma^2 + \|\mu\|^2 \tag{6}$$

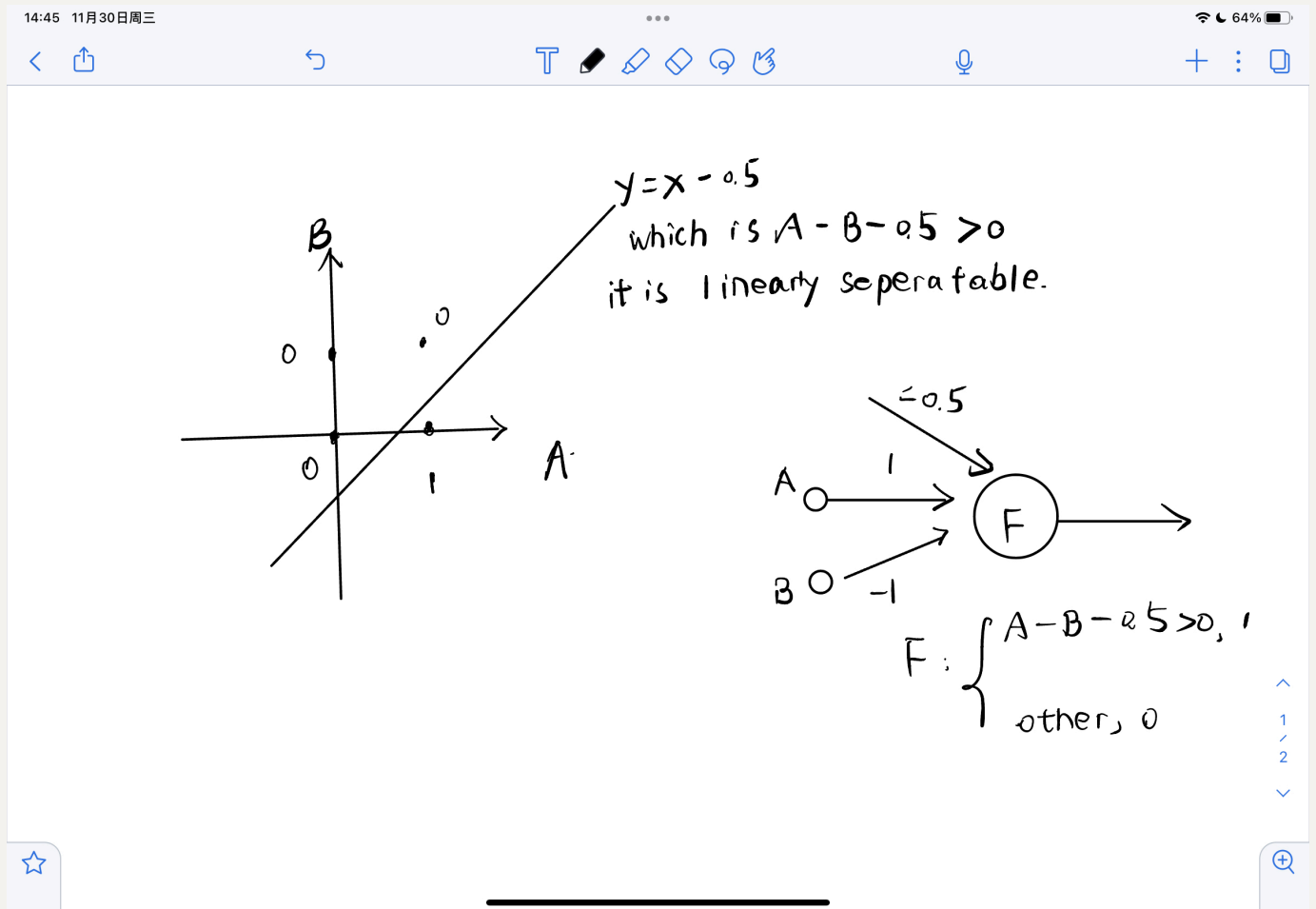
Here L is the dimension of rm we can write.

$$\begin{aligned}
E[t|x] &= \int t p(t|x) dt \\
&= \int t \sum_{k=1}^K \pi_k N(t|\mu_k, \sigma^2) dt \\
&= \sum_{k=1}^K \pi_k \int t N(t|\mu_k, \sigma_k^2) dt \\
&= \sum_{k=1}^K \pi_k \mu_k
\end{aligned} \tag{7}$$

(b)

$$\begin{aligned}
s^2(x) &= E[\|t - E[t|x]\|^2|x] = E[(t^2 - 2tE[t|x] + E[t|x]^2)|x] \\
&= E[t^2|x] - E[2tE[t|x]|x] + E[t|x]^2 = E[t^2|x] - E[t|x]^2 \\
&= \int \|t\|^2 \sum_{k=1}^K \pi_k N(\mu_k, \sigma_k^2) dt - \left\| \sum_{l=1}^K \pi_l \mu_l \right\|^2 \\
&= \sum_{k=1}^K \pi_k (L\sigma_k^2 + \|\mu_k\|^2) - \left\| \sum_{l=1}^K \pi_l \mu_l \right\|^2 \\
&= L \sum_{k=1}^K \pi_k \sigma_k^2 + \sum_{k=1}^K \pi_k \|\mu_k\|^2 \\
&\quad - 2 \left( \sum_{l=1}^K \pi_l \mu_l \right) \left( \sum_{k=1}^K \pi_k \mu_k \right) + \sum_{k=1}^K \pi_k \left\| \sum_{l=1}^K \pi_l \mu_l \right\|^2 \\
&= L \sum_{k=1}^K \pi_k \sigma_k^2 + \sum_{k=1}^K \pi_k \|\mu_k - \sum_{l=1}^K \pi_l \mu_l\|^2 \\
&= \sum_{k=1}^K \pi_k (L\sigma_k^2 + \|\mu_k - \sum_{l=1}^K \pi_l \mu_l\|^2)
\end{aligned} \tag{8}$$

**Q4**



## Q5

(a)  $4 * 4 * 3 = 48$

(b)  $(10/2) * (10/2) * 3 = 75$

(c)  $48 + 5 * 5 * 3 * 4 = 348$

(d) it's right

(e) It needs to learn more parameters.

## Q6

★ **SOLUTION:** Connect the input for  $X_1$  to the output unit with a weight  $C * (w_5 * w_1 + w_6 * w_2)$ , and connect the input for  $X_2$  to the output unit with weight  $C(w_5 * w_3 + w_6 * w_4)$ . Then the output unit can use the same activation function it used originally.

(a) Connect the input for  $X_1$  to the output unit with a weight  $C * (w_5 * w_1 + w_6 * w_2)$ , and connect the input for  $X_2$  to the output unit with weight  $C(w_5 * w_3 + w_6 * w_4)$ . Then the output unit can use the same activation function it used originally.

(b) This is true. Each layer can be thought of as performing a matrix multiply to find its representation given the representation on the layer that it receives input from. Thus the entire network just performs a chain of matrix multiplies, and therefore we can simply multiply the matrices together to find the weights to represent the function with a single layer.

(c) One solution:  $w_1 = w_3 = -10, w_2 = w_4 = -1, w_5 = 5$ , and  $w_6 = -6$ . The intuition here is that we can decompose  $A \text{ XOR } B$  into  $(A \text{ OR } B) \text{ AND NOT } (A \text{ AND } B)$ . We make the upper hidden unit behave like an OR by making it saturate when either of the input units are 1. It isn't possible to make a hidden unit that behaves exactly like AND, but we can at least make the lower hidden unit continue to increase in activation after the upper one has saturated.