

# HW4 Answer

---

## Problem 1

---

<b>t</b>	<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>
$z_i$	90	90	190	90
$f(z_i)$	1	1	1	1
$z_f$	10	10	-90	10
$f(z_i)$	1	1	0	1
$z_o$	-10	90	90	90
$f(z_o)$	0	1	1	1
$z$	3	-2	4	0
$c'$	3	1	4	4
$y$	0	1	4	4

## Problem 2

---

First, give an toy example: Consider a fully connected network with  $\tanh$  as the non-linear activation e.g.

$$\begin{aligned}
 Y &= WX + B \\
 Z &= \tanh(Y) \\
 \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} &= \begin{bmatrix} w_1 & w_2 \\ w_3 & w_4 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} \\
 \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} &= \begin{bmatrix} \tanh(y_1) \\ \tanh(y_2) \end{bmatrix}
 \end{aligned}$$

If  $L$  is the loss of the network and given  $\frac{\partial L}{\partial Z}$  (from the preceding layer), then

$$\begin{aligned}
\frac{\partial L}{\partial Y} &= \begin{bmatrix} \frac{\partial L}{\partial y_1} \\ \frac{\partial L}{\partial y_2} \end{bmatrix} \\
&= \begin{bmatrix} \frac{\partial L}{\partial z_1} (1 - z_1^2) \\ \frac{\partial L}{\partial z_2} (1 - z_2^2) \end{bmatrix} \\
&= \begin{bmatrix} (1 - z_1^2) \\ (1 - z_2^2) \end{bmatrix} \odot \begin{bmatrix} \frac{\partial L}{\partial z_1} \\ \frac{\partial L}{\partial z_2} \end{bmatrix} \\
&= \begin{bmatrix} (1 - z_1^2) & 0 \\ 0 & (1 - z_2^2) \end{bmatrix} \begin{bmatrix} \frac{\partial L}{\partial z_1} \\ \frac{\partial L}{\partial z_2} \end{bmatrix} \\
&= \text{diag}(1 - Z^2) \frac{\partial L}{\partial Z} \\
\frac{\partial L}{\partial W} &= \frac{\partial L}{\partial Y} X^T \\
\frac{\partial L}{\partial X} &= W^T \frac{\partial L}{\partial Y}
\end{aligned}$$

To write explicitly, let  $W_1 = W_2 = W_h$  and  $U_1 = U_2 = W_i$ , then

$$\begin{aligned}
h_1 &= \tanh(W_1 h_0 + U_1 x_1) \\
h_2 &= \tanh(W_2 h_1 + U_2 x_2)
\end{aligned}$$

Let  $z = W_o h_2$ . By the chain rule, we can derive the results as follows:

$$\begin{aligned}
\frac{\partial L}{\partial W_o} &= \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial z}{\partial W_o} \\
\frac{\partial L}{\partial W} &= \frac{\partial L}{\partial W_1} + \frac{\partial L}{\partial W_2} \\
&= \text{diag}(1 - (h_1)^2) \frac{\partial L}{\partial h_1} h_0^T + \text{diag}(1 - (h_2)^2) \frac{\partial L}{\partial h_2} h_1^T \\
\frac{\partial L}{\partial U} &= \frac{\partial L}{\partial U_1} + \frac{\partial L}{\partial U_2} \\
&= \text{diag}(1 - (h_1)^2) \frac{\partial L}{\partial h_1} x_1^T + \text{diag}(1 - (h_2)^2) \frac{\partial L}{\partial h_2} x_2^T
\end{aligned}$$

Note that  $W_i, W_h \in \mathbb{R}^{n \times n}$ ,  $h_2, h_1 \in \mathbb{R}^n$ , and  $W_o \in \mathbb{R}^{1 \times n}$ . Let

$$L(y, \hat{y}) = -y \log \hat{y} - (1 - y) \log (1 - \hat{y})$$

1.  $\frac{\partial L(y, \hat{y})}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z}$ :  
 $\frac{\partial L(y, \hat{y})}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z} = -[y \cdot \frac{1}{\hat{y}} + (1 - y) \cdot \frac{-1}{1 - \hat{y}}] \cdot \hat{y}(1 - \hat{y}) = \hat{y} - y$
2.  $\frac{\partial z}{\partial W_o}$ :  
 $\frac{\partial z}{\partial W_o} = h_2^T$
3.  $\frac{\partial z}{\partial h_2}$ :  $W_o^T$

By (1)(2), we can get  $\frac{\partial L(y, \hat{y})}{\partial W_o} = h_2^T (\hat{y} - y)$

Now,

$$\frac{\partial L}{\partial h_2} = W_o^T \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z}$$

$$\frac{\partial L}{\partial h_1} = W_h^T \text{diag}(1 - (h_2)^2) \frac{\partial L}{\partial h_2}$$

Then, we can get  $\frac{\partial L(y, \hat{y})}{\partial W_h} = (\hat{y} - y) \text{diag}(1 - (h_2)^2) W_o^T h_1^T$ , and

$$\frac{\partial L(y, \hat{y})}{\partial W_i} = (\hat{y} - y) \text{diag}(1 - (h_1)^2) W_h^T \text{diag}(1 - (h_2)^2) W_o^T x_1^T + \text{diag}(1 - (h_2)^2) W_o^T x_2^T.$$

## Problem 3

Given  $\mathbf{g}_{t-1} = \{g_{t-1}^k\}_{k=1}^K$ , we update  $\mathbf{g}_t = \{g_{t-1}^k + \alpha_t f_t^k\}_{k=1}^K = \mathbf{g}_{t-1} + \alpha_t \mathbf{f}_t$  as follows:

$$\begin{aligned} \mathbf{f}_t &\in \underset{f \in \mathcal{F}}{\text{argmin}} \left. \frac{\partial}{\partial \alpha} L(\mathbf{g}_{t-1} + \alpha \mathbf{f}) \right|_{\alpha=0} \\ &= \underset{f \in \mathcal{F}}{\text{argmin}} \frac{\partial}{\partial \alpha} \sum_{i=1}^n \exp \left( \left( \frac{1}{K-1} \sum_{k \neq \hat{y}_i} g_{t-1}^k(x_i) - g_{t-1}^{\hat{y}_i}(x_i) \right) + \alpha \left( \frac{1}{K-1} \sum_{k \neq \hat{y}_i} f^k(x_i) - f^{\hat{y}_i}(x_i) \right) \right) \Big|_{\alpha=0} \\ &= \underset{f \in \mathcal{F}}{\text{argmin}} \sum_{i=1}^n \exp \left( \frac{1}{K-1} \sum_{k \neq \hat{y}_i} g_{t-1}^k(x_i) - g_{t-1}^{\hat{y}_i}(x_i) \right) \left( \frac{1}{K-1} \sum_{k \neq \hat{y}_i} f^k(x_i) - f^{\hat{y}_i}(x_i) \right) \\ &= \underset{f \in \mathcal{F}}{\text{argmin}} Z_t \mathbb{E}_{i \sim D_t} \left[ \frac{1}{K-1} \sum_{k \neq \hat{y}_i} f^k(x_i) - f^{\hat{y}_i}(x_i) \right] = \underset{f \in \mathcal{F}}{\text{argmin}} Z_t \mathbb{E}_{i \sim D_t} \left[ \frac{1}{K-1} \cdot 1\{f(x_i) \neq \hat{y}_i\} - 1\{f(x_i) = \hat{y}_i\} \right] \\ &= \underset{f \in \mathcal{F}}{\text{argmin}} Z_t \left( \frac{K}{K-1} \mathbb{P}_{i \sim D_t} [f(x_i) \neq \hat{y}_i] - 1 \right) = \underset{f \in \mathcal{F}}{\text{argmin}} \mathbb{P}_{i \sim D_t} [f(x_i) \neq \hat{y}_i] \\ \alpha_t &\in \underset{\alpha \in \mathbb{R}}{\text{argmin}} L(\mathbf{g}_{t-1} + \alpha \mathbf{f}_t) \\ &= \underset{\alpha \in \mathbb{R}}{\text{argmin}} \sum_{i=1}^n \exp \left( \left( \frac{1}{K-1} \sum_{k \neq \hat{y}_i} g_{t-1}^k(x_i) - g_{t-1}^{\hat{y}_i}(x_i) \right) + \alpha \left( \frac{1}{K-1} \sum_{k \neq \hat{y}_i} f_t^k(x_i) - f_t^{\hat{y}_i}(x_i) \right) \right) \\ &= \underset{\alpha \in \mathbb{R}}{\text{argmin}} Z_t \mathbb{E}_{i \sim D_t} \left[ e^{\alpha \left( \frac{1}{K-1} \sum_{k \neq \hat{y}_i} f_t^k(x_i) - f_t^{\hat{y}_i}(x_i) \right)} \right] \\ &= \underset{\alpha \in \mathbb{R}}{\text{argmin}} Z_t \mathbb{E}_{i \sim D_t} \left[ e^{\frac{\alpha}{K-1} \cdot 1\{f_t(x_i) \neq \hat{y}_i\}} + e^{-\alpha} \cdot 1\{f_t(x_i) = \hat{y}_i\} \right] \\ &= \underset{\alpha \in \mathbb{R}}{\text{argmin}} Z_t \left( \epsilon_t e^{\frac{\alpha}{K-1}} + e^{-\alpha} (1 - \epsilon_t) \right) = \left\{ \frac{K-1}{K} \log \frac{(K-1)(1 - \epsilon_t)}{\epsilon_t} \right\} \end{aligned}$$

where

$$Z_t = \sum_{i=1}^n \exp \left( \frac{1}{K-1} \sum_{k \neq \hat{y}_i} g_{t-1}^k(x_i) - g_{t-1}^{\hat{y}_i}(x_i) \right)$$

and that  $D_t$  is a probability distribution for  $t = 1, \dots, n$  given by

$$D_t(i) = \frac{1}{Z_t} \exp \left( \frac{1}{K-1} \sum_{k \neq \hat{y}_i} g_{t-1}^k(x_i) - g_{t-1}^{\hat{y}_i}(x_i) \right)$$

and that  $\epsilon_t = \mathbb{P}_{i \sim D_t} [f_t(x_i) \neq \hat{y}_i]$  is the error of  $f_t$  on training sample weighted by the distribution  $D_t$ .

## Problem 4

Follow the lecture note ([https://ntueemlta2022.github.io/slides/week9/W9\\_GMM\\_EM.pdf](https://ntueemlta2022.github.io/slides/week9/W9_GMM_EM.pdf) ([https://ntueemlta2022.github.io/slides/week9/W9\\_GMM\\_EM.pdf](https://ntueemlta2022.github.io/slides/week9/W9_GMM_EM.pdf))). Calculate

$$\begin{aligned} \pi_k^{(t+1)} &= \frac{1}{N} \sum_{i=1}^N \delta_{ik}^{(t)} \\ \boldsymbol{\mu}_k^{(t+1)} &= \frac{\sum_{i=1}^N \delta_{ik}^{(t)} \mathbf{x}_i}{\sum_{i=1}^N \delta_{ik}^{(t)}} \\ \boldsymbol{\Sigma}_k^{(t+1)} &= \frac{\sum_{i=1}^N \delta_{ik}^{(t)} \left( \mathbf{x}_i - \boldsymbol{\mu}_k^{(t+1)} \right) \left( \mathbf{x}_i - \boldsymbol{\mu}_k^{(t+1)} \right)^T}{\sum_{i=1}^N \delta_{ik}^{(t)}} \end{aligned}$$

explicitly to get all points.

The calculation process follows the HW2 answer. (之後“可能”會補上計算過程)