

#3 $\ell^{\text{CE}}(f, \gamma) = -\log \left(\frac{\exp(f_\gamma)}{\sum_{j=1}^k \exp(f_j)} \right)$, $f \in \mathbb{R}^k$, $\gamma \in \{1, \dots, k\}$

(a) $\frac{\exp(f_\gamma)}{\sum_{j=1}^k \exp(f_j)} = \frac{\exp(f_\gamma)}{\exp(f_1) + \exp(f_\gamma) + \exp(f_k)}$ Since $\exp(f_i) > 0$ for all $i \in \{1, \dots, k\}$,
 $0 < \frac{\exp(f_\gamma)}{\sum \exp(f_j)} < 1$

Also, we know that $-\infty < \log(x) < 0$ when $x \in (0, 1)$

Therefore, $-\infty < \log \left(\frac{\exp(f_\gamma)}{\sum \exp(f_j)} \right) < 0$ and thus, $0 < \ell^{\text{CE}}(f, \gamma) < \infty$

(b) $\ell^{\text{CE}}(\lambda e_\gamma, \gamma) = -\log \left(\frac{\exp(\lambda)}{(k-1) + \exp(\lambda)} \right) = \log \left(\frac{e^\lambda + k-1}{e^\lambda} \right) = \log \left(1 + \frac{k-1}{e^\lambda} \right) \xrightarrow{\lambda \rightarrow \infty} \log(1) = 0!$

#4 Let $f(x) = \max \{f_1(x), \dots, f_k(x)\}$, $f_i(x)$: diff, univariate fh.

pf) for a given x , maximizing index $I = \arg \max_i \{f_i(x)\}$ is unique.

Therefore, $f_I(x) > f_i(x)$ holds for all $i \neq I$

Since $\forall i$, $f_i(x)$ is differentiable, $\forall i$ $f_i(x)$ is continuous. $\therefore \forall i \neq I$ $f_I(x) - f_i(x)$ is continuous.

By the definition of continuity, $f_I(x) - f_i(x) > 0$, $\exists \delta > 0$ s.t. if $x - \delta < t < x + \delta$, $f_I(t) - f_i(t) > 0$ holds.

This means that $\forall t \in (x - \delta, x + \delta)$, $f_I(t) = f(t)$.

Therefore for $|h| < \delta$, $\lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} = \lim_{h \rightarrow 0} \frac{f_I(x+h) - f_I(x)}{h}$ holds.

i.e. $f'(x) = f_I'(x)$!

#5 (a) $r(z) = \max\{0, z\}$, $r'(z) = r(z)$

1) $r'(z) = \max\{0, r(z)\} = \max\{0, \max\{0, z\}\} = \max\{0, 0, z\} = \max\{0, z\} = r(z)$

(b) Softplus $r(z) = \log(1+e^z)$, $r'(z) = \frac{e^z}{1+e^z}$

Let's show that $r(z)$ has Lipschitz continuous derivative

$$|r'(x) - r'(y)| = \left| \frac{e^x}{1+e^x} - \frac{e^y}{1+e^y} \right| = \left| \frac{e^x + e^x e^y - e^y - e^y e^x}{(1+e^x)(1+e^y)} \right| = \left| \frac{e^x - e^y}{(1+e^x)(1+e^y)} \right|$$

$$= \left| \frac{e^c}{(1+e^x)(1+e^y)} \right| |x-y| \quad \exists c, c \in (\min(x, y), \max(x, y)) \text{ by Mean Value Theorem}$$

Since $c < \max(x, y)$, $e^c < 1+e^x$ or $e^c < 1+e^y$ holds,

$\therefore |r'(x) - r'(y)| = \left| \frac{e^c}{(1+e^x)(1+e^y)} \right| |x-y| < |x-y|$ \therefore Softplus has Lipschitz continuous derivative.

on the other hand, Consider ReLU, $r(z) = \max(0, z)$

$$\frac{|r'(z) - r'(-z)|}{|z - (-z)|} = \frac{+1}{2|z|} \rightarrow \infty \text{ as } z \rightarrow 0$$

$\therefore \forall L, \exists \delta > 0$ s.t if $0 < |z| < \delta$, $|r'(z) - r'(-z)| > L|z - (-z)|$. \therefore ReLU does not have Lipschitz continuous derivative!

(c) $p(z) = \frac{1-e^{-2z}}{1+e^{-2z}} = \frac{e^{2z}-1}{e^{2z}+1} = \frac{2e^{2z} - (1+e^{2z})}{1+e^{2z}} = 2\frac{1}{1+e^{-2z}} - 1 = 2r(2z) - 1$

$L > 1$, $A_1, \dots, A_L, b_1, \dots, b_L$ are given, $y_L = A_L y_{L-1} + b_L, \dots, y_1 = r(A_1 x + b_1)$ 이라 하자.

① $y_L = A_L r(A_{L-1} y_{L-2} + b_{L-1}) + b_L = C_L r(C_{L-1} y_{L-2} + d_{L-1}) + d_L$

$\Rightarrow \boxed{C_L = \frac{1}{2} A_L, d_L = b_L + \frac{1}{2} A_L \begin{pmatrix} 1 \\ 1 \end{pmatrix}}$ $= C_L \left(2r\left(\boxed{2C_{L-1} y_{L-2} + 2d_{L-1}}\right) \right) + d_L - C_L \begin{pmatrix} 1 \\ 1 \end{pmatrix}$

$A_{L-1} r(A_{L-2} y_{L-3} + b_{L-2}) + b_{L-1} = 2C_{L-1} + p(C_{L-2} y_{L-3} + d_{L-2}) + 2d_{L-1}$

$= 2C_{L-1} \left(2r(2C_{L-2} y_{L-3} + 2d_{L-2}) - \begin{pmatrix} 1 \\ 1 \end{pmatrix} \right) + 2d_{L-1}$

$\Rightarrow C_{L-1} = \frac{1}{4} A_{L-1}, d_{L-1} = \frac{1}{2} \left(b_{L-1} + \frac{1}{2} A_{L-1} \begin{pmatrix} 1 \\ 1 \end{pmatrix} \right)$ $< C_{L-1} r\left(\boxed{2C_{L-2} y_{L-3} + 2d_{L-2}}\right) + 2d_{L-1} - 2C_{L-1} \begin{pmatrix} 1 \\ 1 \end{pmatrix}$

$= \frac{1}{2} b_{L-1} + \frac{1}{4} A_{L-1} \begin{pmatrix} 1 \\ 1 \end{pmatrix}$ \downarrow 같은 꼴이 계속!

② $\boxed{C_i = \frac{1}{2^i} A_i, d_i = \frac{1}{2} b_i + \frac{1}{2^i} A_i \begin{pmatrix} 1 \\ 1 \end{pmatrix}} \quad (1 \leq i \leq L)$ $\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix} \text{ is a fixed } \mathbb{R}^1 \text{ vector} \right)$

③ $\boxed{C_1 = \frac{1}{2} A_1, d_1 = \frac{1}{2} b_1}$

여기 A_i, b_i 에 대해서, 각각 같은 C_i, d_i 를 사용하면 등가입니다!

pf) $\langle i \rangle L=2$ 인 경우

$$y_2 = A_2 y_1 + b_2$$

$$y_1 = \sigma(A_1 x + b_1)$$

$$\bar{y}_2 = C_2 \bar{y}_1 + d_2 = \frac{1}{2} A_2 \bar{y}_1 + b_2 + \frac{1}{2} A_2 \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

$$\bar{y}_1 = \rho(C_1 x + d_1) = \rho\left(\frac{1}{2} A_1 + \frac{1}{2} b_1\right)$$

$$\begin{aligned} \bar{y}_2 &= \frac{1}{2} A_2 \left(2 \sigma\left(\frac{1}{2} A_1 x + \frac{1}{2} b_1 \right) - \begin{pmatrix} 1 \\ 1 \end{pmatrix} \right) + b_2 + \frac{1}{2} A_2 \begin{pmatrix} 1 \\ 1 \end{pmatrix} \\ &= A_2 \sigma(A_1 x + b_1) + b_2 = y_2. \end{aligned}$$

$\langle ii \rangle L=3$ 인 경우

$$y_3 = A_3 y_2 + b_3 \quad \bar{y}_3 = \frac{1}{2} A_3 \bar{y}_2 + b_3 + \frac{1}{2} A_3 \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

$$y_2 = \sigma(A_2 y_1 + b_2) \quad \bar{y}_2 = \rho\left(\frac{1}{4} A_2 \bar{y}_1 + \frac{1}{2} b_2 + \frac{1}{4} A_2 \begin{pmatrix} 1 \\ 1 \end{pmatrix}\right)$$

$$y_1 = \sigma(A_1 x + b_1) \quad \bar{y}_1 = \rho\left(\frac{1}{2} A_1 x + \frac{1}{2} b_1\right) = 2 \sigma(A_1 x + b_1) - \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

$$\bar{y}_2 = \rho\left(\frac{2}{4} A_2 \sigma(A_1 x + b_1) - \frac{1}{4} A_2 \begin{pmatrix} 1 \\ 1 \end{pmatrix} + \frac{1}{2} b_2 + \frac{1}{4} A_2 \begin{pmatrix} 1 \\ 1 \end{pmatrix}\right)$$

$$= 2 \sigma\left(A_2 \sigma(A_1 x + b_1) + b_2\right) - \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

$$\bar{y}_3 = A_3 \sigma\left(A_2 \sigma(A_1 x + b_1) + b_2\right) + b_3 = y_3!$$

$\langle iii \rangle L$ 일때 성립 가정

$$\therefore \underline{y_L} = A_L y_{L-1} + b_L = \underline{\bar{y}_L} = \frac{1}{2} A_L \bar{y}_{L-1} + b_L + \frac{1}{2} A_L \begin{pmatrix} 1 \\ 1 \end{pmatrix} \text{ 인 상반. | 같 하리. } \dots (*)$$

$C_L = \frac{1}{2} A_L$ 이고, $d_L = \frac{1}{2} b_L + \frac{1}{4} A_L \begin{pmatrix} 1 \\ 1 \end{pmatrix}$ 이 바꿀때 $y_{L-1} = \bar{y}_{L-1}$ 임을 보임으로서 $L+1$ 일때 성립함을 보일 수 있다.

$$\begin{aligned} A_{L+1} \sigma(A_L y_{L-1} + b_L) + b_{L+1} &= \frac{1}{2} A_{L+1} \rho\left(\frac{1}{4} A_L \bar{y}_{L-1} + \frac{1}{2} b_L + \frac{1}{4} A_L \begin{pmatrix} 1 \\ 1 \end{pmatrix}\right) + b_{L+1} + \frac{1}{2} A_{L+1} \begin{pmatrix} 1 \\ 1 \end{pmatrix} \\ &= A_{L+1} \sigma\left(\frac{1}{2} A_L \bar{y}_{L-1} + b_L + \frac{1}{2} A_L \begin{pmatrix} 1 \\ 1 \end{pmatrix}\right) + b_{L+1} \end{aligned}$$

$$\Leftrightarrow A_L y_{L-1} + b_L = \frac{1}{2} A_L \bar{y}_{L-1} + b_L + \frac{1}{2} A_L \begin{pmatrix} 1 \\ 1 \end{pmatrix} \dots (*) \text{ 와 동일한 식! } \therefore \text{ 성립한다!}$$

즉 $\langle i \rangle, \langle ii \rangle, \langle iii \rangle, M.I$ 에 대하여 \square page 에 주어진대로 C_i, d_i 를 참으로 동일한 수열을 얻게 된다는 사실을 간단히 증명하였다.

#6. $a_j^0 x_i + b_j^0 < 0 \quad \forall i \in \{1, \dots, N\}. \quad \sigma = \text{ReLU}$

We want to show that $a_j^k x_i + b_j^k < 0$ for all i and k (to show that j th ReLU output remains dead)

$$\begin{aligned} \frac{\partial}{\partial a_j} l(f_\theta(x_i), y_i) &= \frac{\partial l}{\partial f} \frac{\partial f}{\partial a_j} = \frac{\partial l}{\partial f} \frac{\partial u^T v(a x_i + b)}{\partial a_j} = \frac{\partial l}{\partial f} u_j \cdot v'(a_j x_i + b_j) x_i \\ &= \begin{cases} 0 & \text{if } a_j x_i + b_j < 0 \\ \frac{\partial l}{\partial f} u_j x_i & \text{if } a_j x_i + b_j > 0 \end{cases} \end{aligned}$$

$$\begin{aligned} \frac{\partial}{\partial b_j} l(f_\theta(x_i), y_i) &= \frac{\partial l}{\partial f} \frac{\partial f}{\partial b_j} = \frac{\partial l}{\partial f} \frac{\partial u^T v(a x_i + b)}{\partial b_j} = \frac{\partial l}{\partial f} u_j v'(a_j x_i + b_j) \\ &= \begin{cases} 0 & \text{if } a_j x_i + b_j < 0 \\ \frac{\partial l}{\partial f} u_j & \text{if } a_j x_i + b_j > 0 \end{cases} \end{aligned}$$

If we think about SGD, $\theta_{k+1} = \theta_k - \text{lr} \times \frac{\partial l}{\partial \theta_k}$.

Since $\langle i \rangle \quad a_j^0 x_i + b_j^0 < 0 \quad (\forall i)$,

$\langle ii \rangle \quad a_j^k x_i + b_j^k < 0$. Then, $\frac{\partial l}{\partial a_j} = 0$ and $\frac{\partial l}{\partial b_j} = 0$ for all i .

Thus, a_j^{k+1}, b_j^{k+1} remains unchanged and thus $a_j^{k+1} x_i + b_j^{k+1} < 0$ for all i .

By Mathematical induction, $\langle i \rangle, \langle ii \rangle$, we proved that $\forall i, k, a_j^k x_i + b_j^k < 0$ holds.

The j th ReLU output is dead throughout training.

#7 If we use Leaky ReLU, $\sigma(z) = \begin{cases} z & z \geq 0 \\ \alpha z & z < 0 \end{cases}$

$$\frac{\partial \ell(f_{\theta}(x_i), y_i)}{\partial a_j} = \begin{cases} \left(\frac{\partial \ell}{\partial f} u_j x_i \right) \alpha & (a_j x_i + b_j < 0) \\ \left(\frac{\partial \ell}{\partial f} u_j x_i \right) & (a_j x_i + b_j \geq 0) \end{cases} \quad \frac{\partial \ell(f_{\theta}(x_i), y_i)}{\partial b_j} = \begin{cases} \left(\frac{\partial \ell}{\partial f} u_j \right) \alpha & (a_j x_i + b_j < 0) \\ \frac{\partial \ell}{\partial f} u_j & (a_j x_i + b_j \geq 0) \end{cases}$$

$\alpha \neq 0$ 이라면, $\forall i, a_j x_i + b_j < 0$ 일때

$$\frac{\partial \ell}{\partial a_j} = \left(\frac{\partial \ell}{\partial f} u_j x_i \right) \alpha : \text{not identically zero}$$

$$\frac{\partial \ell}{\partial b_j} = \left(\frac{\partial \ell}{\partial f} u_j \right) \alpha : \text{not identically zero} \quad \text{이때}$$

SGD 에 따라 $a_j^{k+1} = a_j^k - \frac{\partial \ell}{\partial a_j} \times \text{learning_rate}$

$$b_j^{k+1} = b_j^k - \frac{\partial \ell}{\partial b_j} \times \text{learning_rate} \quad \text{이 때에 update 되며}$$

#6 보 문제의 경우 gradient 가 "exactly vanish" 하지는 않는다!