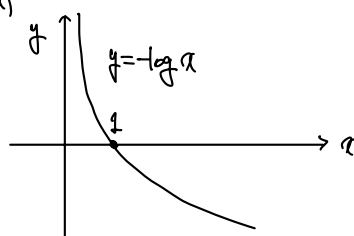


#3.

(a)

This is the graph of $f(x) = -\log x$ $\forall x \in (0, 1)$, $0 < f(x) < \infty$ since $f(x) = -\frac{1}{x} < 0$ f is strictly decreasing, and $f(1) = 0$. $\therefore f(x) > 0$ for $x \in (0, 1)$.

$$\text{Since } 0 < \frac{\exp(f_1)}{\sum_{j=1}^k \exp(f_j)} < 1, \quad 0 < l^{\text{CE}}(f, y) < \infty.$$

■

$$(b) \quad l^{\text{CE}}(\lambda e_y, y) = -\log \left(\frac{\exp(\lambda)}{\exp(\lambda)} \right) = 0 \quad \text{for } \forall \lambda.$$

$$\therefore l^{\text{CE}}(\lambda e_y, y) \rightarrow 0 \quad \text{as } \lambda \rightarrow \infty.$$

■

#4.

Since I is unique at $x \in \mathbb{R}$, $\exists \delta > 0$ s.t. $\forall t \in (x-\delta, x+\delta)$, $\arg \max_{i \in \{1, N\}} \{f_i(t)\} = I$.

(pf) let's assume that $\forall \delta > 0$, $\exists t \in (x-\delta, x+\delta)$ s.t. $\arg \max_{i \in \{1, N\}} \{f_i(t)\} = J \neq I$

$f_I(t) < f_J(t)$, and take $\delta \rightarrow 0$. Then $t \rightarrow x$, and since f_I and f_J are continuous, $\lim_{t \rightarrow x} f_I(t) \leq \lim_{t \rightarrow x} f_J(t)$ means $f_I(x) \leq f_J(x)$ which is contradiction to I being the argmax of $\{f_i(x)\}$ i.e. $f_I(x) > f_J(x)$.

■

$$\text{Therefore, } \lim_{h \rightarrow 0} \frac{f_I(x+h) - f_I(x)}{h} = \lim_{h \rightarrow 0} \frac{f_I(x+h) - f_J(x)}{h} = f'_I(x).$$

(and $h < \delta$).

■

#5.

$$(a) \quad (1) \quad z \geq 0: \quad \sigma(\sigma(z)) = \sigma(z)$$

$$(2) \quad z < 0: \quad \sigma(\sigma(z)) = \sigma(0) = 0 = \sigma(z)$$

■

$$(b) \quad \sigma'(z) = \frac{e^z}{1+e^z}, \quad \sigma''(z) = \frac{e^z}{(1+e^z)^2} \leq 1$$

$$\forall x, y \in \mathbb{R}, \quad |\sigma'(x) - \sigma'(y)| = |\sigma''(\xi)| |x-y| \leq |x-y| \Rightarrow \text{Lip schitz continuous.}$$

\uparrow
($\xi \in (x, y)$, by MVT)

ReLu: first, the derivative is not well defined, but if we define it as follows

$$\text{ReLu}'(x) = \begin{cases} 0 & (x < 0) \\ 1 & (x \geq 0) \end{cases}, \quad \text{this isn't Lipschitz since if we assume it is for } x > 0,$$

$$\text{Let } \varepsilon = \left[\frac{1}{2K}, \frac{1}{2K} \right], \text{ and take } x_1 = -\varepsilon, \quad x_2 = \varepsilon$$

$$\rightarrow |f(x_1) - f(x_2)| = 1 \leq K \cdot 2\varepsilon < 1 \rightarrow \text{contradiction, not Lipschitz}$$

■

(c) \Leftarrow
 Given $A_1, \dots, A_L, b_1, \dots, b_L$, $y_1 \sim y_L$ are the states of the MLP with sigmoid.

$$\sigma(z) = \frac{1}{2} \rho\left(\frac{z}{2}\right) + \frac{1}{2}$$

Set $C_1 = \frac{1}{2}A_1, d_1 = \frac{1}{2}b_1$, then $y'_1 := \rho(C_1 x + d_1) = 2y_1 - v_1$ ($v_1 = (1, 1, \dots, 1)^T \in \mathbb{R}^{n_1}$)

Set $C_2 = \frac{1}{4}A_2, d_2 = \frac{1}{4}A_2 v_1 + \frac{1}{2}b_2$, then $y'_2 := \rho(C_2 y'_1 + d_2) = 2y_2 - v_2$ ($v_2 = (1, \dots, 1)^T \in \mathbb{R}^{n_2}$)

\vdots $C_i = \frac{1}{4}A_i, d_i = \frac{1}{4}A_i v_{i-1} + \frac{1}{2}b_i, (2 \leq i \leq L-1)$ ($v_i \in (1, \dots, 1)^T \in \mathbb{R}^{n_i}$)

then $y'_{L-1} := \rho(C_{L-1} y'_{L-2} + d_{L-1}) = 2y_{L-1} - v_{L-1}$.

Set $C_L = \frac{1}{2}A_L, d_L = \frac{1}{2}A_L v_{L-1} + b_L$

$\Rightarrow y'_L := \frac{1}{2}A_L (2y_{L-1} - v_{L-1}) + \frac{1}{2}A_L v_{L-1} + b_L = A_L y_{L-1} + b_L = y_L \Rightarrow$ this new MLP with tanh represent the given MLP

\Leftarrow

Given $C_1 \sim C_L, d_1 \sim d_L$, and states $y_1 \sim y_L$ with MLP with tanh,

(v_i 's are defined same as above)

$$\rho(z) = 2\sigma(2z) - 1$$

Set $A_1 = 2C_1, b_1 = 2d_1$, then $y'_1 := \sigma(A_1 x + b_1) = \frac{1}{2}(y_1 + v_1)$.

Set $A_2 = 4C_2, b_2 = 2d_2 - 2C_2 v_1$, then $y'_2 := \sigma(A_2 y'_1 + b_2) = \frac{1}{2}(y_2 + v_2)$

\vdots Set $A_i = 4C_i, b_i = 2d_i - 2C_i v_{i-1}, (2 \leq i \leq L-1)$

then $y'_{L-1} := \frac{1}{2}(y_{L-1} + v_{L-1})$

Set $A_L = 2C_L, b_L = d_L - C_L v_{L-1}$, then $y'_L := C_L y'_{L-1} + d_L = y_L \Rightarrow$ this new MLP with sigmoid represent the given MLP

#6 In the initialization step, if $a_j^0 x_i + b_j^0 < 0$ for $\forall i \in [1, N]$

$$\frac{\partial}{\partial a_j} \ell(f_0(x_i), y_i) = \frac{\partial}{\partial f_0} \ell(f_0(x_i), y_i) \cdot \frac{\partial f_0}{\partial a_j}(x_i) = \frac{\partial}{\partial f_0} \ell(f_0(x_i), y_i) \cdot u_j \sigma'(a_j x_i + b_j) \cdot x_i$$

$$\frac{\partial}{\partial b_j} \ell(f_0(x_i), y_i) = \frac{\partial}{\partial f_0} \ell(f_0(x_i), y_i) \cdot u_j \sigma'(a_j x_i + b_j)$$

Consider the k^{th} iteration of the SGD.

when $k=1 \Rightarrow$ since $a_j^0 x_i + b_j^0 < 0$, $\sigma'(a_j^0 x_i + b_j^0) = 0$, and $\frac{\partial}{\partial a_j} \ell(f_0(x_i), y_i) = \frac{\partial}{\partial b_j} \ell(f_0(x_i), y_i) = 0$

$\Rightarrow a_j^1 x_i + b_j^1 < 0$.

Assume $a_j^k x_i + b_j^k < 0 \Rightarrow a_j^{k+1} = a_j^k, b_j^{k+1} = b_j^k$ since the gradients are 0 vectors.

\therefore By mathematical induction, $\forall k, a_j^k x_i + b_j^k < 0$, hence the gradient vanishes.

#7.

$$\frac{\partial}{\partial a_j} \ell(f_0(x_i), y_i) = \frac{\partial}{\partial f_0} \ell(f_0(x_i), y_i) \cdot u_j \sigma'(a_j x_i + b_j) \cdot x_i$$

$$\frac{\partial}{\partial b_j} \ell(f_0(x_i), y_i) = \frac{\partial}{\partial f_0} \ell(f_0(x_i), y_i) \cdot u_j \sigma'(a_j x_i + b_j)$$

But, $\sigma'(a_j^0 x_i + b_j^0) = \alpha$, and $\frac{\partial}{\partial a_j} \ell(f_0(x_i), y_i) = \frac{\partial}{\partial f_0} \ell(f_0(x_i), y_i) \cdot (\alpha u_j x_i) \neq 0$, and similar for $\frac{\partial \ell}{\partial b_j}$.

Therefore, there are SGD updates to the parameters a_j and b_j , and the gradient no longer vanishes. ■