

3.

$$(a) \quad 0 < \exp(f_0) < \sum_{j=1}^k \exp(f_j) < \infty$$

$$\Rightarrow 0 < \frac{\exp(f_0)}{\sum_{j=1}^k \exp(f_j)} < 1$$

$$\Rightarrow 0 < -\log\left(\frac{\exp(f_0)}{\sum_{j=1}^k \exp(f_j)}\right) < \infty$$

$$(b) \quad l^{CE}(\lambda e_y, y) = -\log\left(\frac{\exp(\lambda)}{\exp(\lambda) + (k-1)\exp(0)}\right)$$

$$= -\log \frac{e^\lambda}{e^\lambda + (k-1)} \xrightarrow{\text{as } \lambda \rightarrow \infty} 0$$

$$\text{Since } \frac{e^\lambda}{e^\lambda + (k-1)} \xrightarrow{\text{as } \lambda \rightarrow \infty} 1$$

□

4. Let $x \in \mathbb{R}$ be given and $I = \arg\max_i \{f_i(x)\}$ is unique.

Since \mathbb{R} is compact,

$\{y \in \mathbb{R} : f_I(y) > f_n(y), n = \{1, 2, \dots, k\} \setminus \{I\}\}$ is open.

$\therefore \exists \delta > 0$ s.t.

$$\forall y \in \mathbb{R}, \|x-y\| < \delta \Rightarrow f_I(y) > f_n(y), n = \{1, 2, \dots, k\} \setminus \{I\}$$

$$\Leftrightarrow f(y) = f_I(y)$$

$$\therefore \forall y \in \mathbb{R}, \|x-y\| < \delta \Rightarrow f'(y) = f'_I(y)$$

Hence, $f'(x) = f'_I(x)$

□

5.

(a) i) The case $z \geq 0$

$$\sigma(\sigma(z)) = \sigma(z) = z$$

ii) The case $z < 0$

$$\sigma(\sigma(z)) = \sigma(0) = 0 = \sigma(z)$$

$$\therefore \forall z \in \mathbb{R}, \quad \sigma(\sigma(z)) = \sigma(z)$$

(b) $\sigma(z) = \log(1 + e^z)$

$$\sigma'(z) = \frac{e^z}{1 + e^z}$$

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

 $\forall x, y \in \mathbb{R} \exists t \in (x, y) \text{ s.t.}$ $x \neq y$

$$\frac{\sigma'(x) - \sigma'(y)}{x - y} = \sigma''(t) = \frac{e^t}{(1 + e^t)^2} < 1$$

$$\therefore |\sigma'(x) - \sigma'(y)| < |x - y|$$

 $\therefore \sigma(z)$ has Lipschitz continuous derivatives.

On the other hand, ReLU does not

Since ReLU is non-differentiable at $z = 0$.(c) $p(z) = 2\sigma(2z) - \underline{1}$

$$\therefore y_1 = p(c_1 z + d_1) = 2\sigma(\underbrace{2c_1 z}_{c'_1} + d_1) - \underline{1} = 2\sigma(c'_1 z + d_1) - \underline{1}$$

$$y_2 = p(c_2 y_1 + d_2) = 2\sigma(\underbrace{2c_2}_{\text{!!}} \underbrace{\sigma(c'_1 z + d_1)}_{y'_1} + \underbrace{(d_2 - c_2 \underline{1})}_{\text{!!}}) - \underline{1}$$

$$= 2\sigma(c'_2 y'_1 + d_2) - \underline{1}$$

⋮

$$y_{L-1} = 2\sigma(c_{L-1}' y_{L-2}' + d_{L-1}') - \underline{1}$$

$$y_L = c_L y_{L-1} + d_L = 2c_L \sigma(c_{L-1}' y_{L-2}' + d_{L-1}') + d_L - c_L \underline{1} = c_L' \sigma(y_{L-1}') + d_L'$$

By initializing $\begin{cases} c'_l = 2c_l = A_l \in \mathbb{R}^{n_l \times n_{l-1}} \\ d'_l = d_l - c_l \underline{1} = b_l \in \mathbb{R}^{n_l} \end{cases} \quad (l=1, \dots, L)$,

y'_l be equivalent with y_l in MLPs built with Sigmoid activations.
($l=1, \dots, L$)



$$6. \quad x \in \mathbb{R} \setminus \{0\} \Rightarrow \begin{cases} \nabla_b f_\theta(x) = \sigma'(ax+b) \odot u \\ \nabla_a f_\theta(x) = (\sigma'(ax+b) \odot u) \cdot x \end{cases}$$

$$a_j^k x_i + b_j^k < 0 \Rightarrow \sigma'(a_j^k x_i + b_j^k) = 0$$

$$\Rightarrow \begin{cases} \nabla_{b_j} f_\theta(x) = 0 \\ \nabla_{a_j} f_\theta(x) = 0 \end{cases}$$

$$\Rightarrow \begin{cases} \nabla_{b_j} l(f_\theta(x), y_i) = \nabla_{b_j} f_\theta(x) \nabla_{f_\theta(x)} l(f_\theta(x), y_i) = 0 \\ \nabla_{a_j} l(f_\theta(x), y_i) = \nabla_{a_j} f_\theta(x) \nabla_{f_\theta(x)} l(f_\theta(x), y_i) = 0 \end{cases}$$

\therefore If $a_j^0 x_i + b_j^0 < 0$ for all i at initialization,

then a_j^k, b_j^k will never update with SGD and

$a_j^k x_i + b_j^k$ remains negative for all $k=1, 2, \dots ; i$,

which result in $\sigma(a_j^k x_i + b_j^k) = 0$ throughout the training.

□

$$7. \quad x \in \mathbb{R} \setminus \{0\} \Rightarrow \left\{ \begin{array}{l} \nabla_a f_\theta(x) = \sigma(ax+b) \\ \nabla_b f_\theta(x) = \sigma'(ax+b) \odot u \\ \nabla_u f_\theta(x) = (\sigma'(ax+b) \odot u) x \end{array} \right.$$

$$a_j^k x_i + b_j^k < 0 \Rightarrow \left\{ \begin{array}{l} \sigma(a_j^k x_i + b_j^k) = 0.01(a_j^k x_i + b_j^k) \\ \sigma'(a_j^k x_i + b_j^k) = 0.01 \end{array} \right.$$

$$\Rightarrow \left\{ \begin{array}{l} \nabla_{a_j} f_\theta(x) = 0.01(a_j^k x_i + b_j^k) \\ \nabla_{b_j} f_\theta(x) = 0.01 u_j \\ \nabla_u f_\theta(x) = 0.01 u_j x \end{array} \right.$$

$$\Rightarrow \nabla_{b_j} l(f_\theta(x), y_i) = \nabla_{b_j} f_\theta(x) \nabla_{f_\theta(x)} l(f_\theta(x), y_i)$$

$$\nabla_{a_j} l(f_\theta(x), y_i) = \nabla_{a_j} f_\theta(x) \nabla_{f_\theta(x)} l(f_\theta(x), y_i)$$

both are not identically zero

with proper loss function.

$\therefore a_j, b_j$ can be updated with SGD,

hence $\sigma(a_j^k x_i + b_j^k)$ can be nonzero with some i, k .

