

MFNN

Homework 3

24/03/21

3. (a) ① Does  $-\log\left(\frac{\exp(f_y)}{\sum_{j=1}^k \exp(f_j)}\right) > 0$  hold?

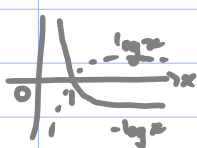
$$\Leftrightarrow \frac{\sum_{j=1}^k \exp(f_j)}{\exp(f_y)} > 1 \quad \text{by applying exp to both sides (exp is increasing } \Rightarrow \text{ preserves inequality)}$$

$$\Leftrightarrow \sum_{j=1}^k \exp(f_j) > \exp(f_y)$$

$$\Leftrightarrow \sum_{\substack{j \in \{1, \dots, k\} \setminus \{y\} \\ =: K'}} \exp(f_j) > 0$$

$\exp(x)$  is strictly positive for all  $x \in \mathbb{R}$  so sum is also strictly positive and hence original exp. holds, i.e.  $\ell^{\text{CE}}(f, y) > 0$

② We need to show  $-\log\left(\frac{\exp(f_y)}{\sum_{j=1}^k \exp(f_j)}\right) < \infty$ , i.e. is finite/bounded above  $\forall f, y, k$



$$-\log\left(\frac{\exp(f_y)}{\sum_{j=1}^k \exp(f_j)}\right) = \log\left(\frac{\sum_{j=1}^k \exp(f_j)}{\exp(f_y)}\right)$$

$$= \log\left(\sum_{j=1}^k \exp(f_j)\right) - f_y \leq \log(k \exp(f_M)) - f_y$$

with  $M := \arg \max_{i \in \{1, \dots, k\}} f_i$

i.e.  $M$  st.  $f_M$  is largest component of  $f$

$$= \log k + f_M - f_y < \infty \quad \text{i.e. is finite since } k, f_M, f_y \in \mathbb{R}$$

$$\text{So } 0 < \ell^{\text{CE}}(f, y) < \infty \quad \blacksquare$$

(b)  $\ell^{\text{CE}}(\lambda \underline{e}_y, y) = -\log\left(\frac{\exp((\lambda \underline{e}_y)_y)}{\sum_{j=1}^k \exp((\lambda \underline{e}_y)_j)}\right)$

$$= -\log\left(\frac{\exp(\lambda)}{\exp(\lambda) + \sum_{j \in K'} \exp(0)}\right) = -\log\left(\frac{\exp(\lambda)}{\exp(\lambda) + (k-1)}\right)$$

$$\begin{aligned} -\log \text{ is continuous so } \lim_{\lambda \rightarrow \infty} -\log(\dots) &= -\log\left(\lim_{\lambda \rightarrow \infty} \left(\frac{\exp \lambda}{\exp \lambda + k-1}\right)\right) \\ &= -\log\left(\lim_{\lambda \rightarrow \infty} \left(\frac{1}{1 + (k-1)e^{-\lambda}}\right)\right) \\ &= -\log\left(\frac{1}{1 + (k-1)(0)}\right) \\ &= -\log(1) = 0 \end{aligned}$$

$$\text{So } \ell^{\text{CE}}(\lambda \underline{e}_y, y) \rightarrow 0 \text{ as } \lambda \rightarrow \infty \quad \blacksquare$$

↙ f differentiable if this limit exists

4.  $\frac{d}{dx} f(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} = \lim_{h \rightarrow 0} \frac{f(x+h) - f_I(x)}{h}$  using given assumption on  $f$  at given  $x \in \mathbb{R}$

If  $I$  is unique at  $x \in \mathbb{R} \Rightarrow$  there is a small neighbourhood,  $(x-\delta, x+\delta)$ ,  $\delta > 0$ , where  $I$  remains unique.  
i.e. a neighbourhood where  $f(x) = f_I(x)$  exactly.

(If there was no  $\delta > 0$  for which this was true, then necessarily another function  $f_i$  with the value of  $f_I$  exists at  $x$  and then  $I$  wouldn't be unique - contradiction)

So for all  $h$  st.  $0 < |h| < \delta$ ,  $f(x+h) = f_I(x+h)$ , hence

$$= \lim_{h \rightarrow 0} \frac{f(x+h) - f_I(x)}{h} = \lim_{h \rightarrow 0} \frac{f_I(x+h) - f_I(x)}{h} = \frac{d}{dx} f_I(x) \quad \blacksquare$$

↑ since  $f_i$  are all differentiable, differential exists for  $i=I$

5.(a) Given  $z \in \mathbb{R}$ ,  $\sigma(z) := \max\{0, z\}$ ,

If  $z > 0$ ,  $\sigma(z) = z \Rightarrow \sigma(\sigma(z)) = \sigma(z) = z = \sigma(z)$

If  $z \leq 0$ ,  $\sigma(z) = 0 \Rightarrow \sigma(\sigma(z)) = \sigma(0) = 0 = \sigma(z) \Rightarrow \sigma(\sigma(z)) = \sigma(z) \quad \blacksquare$

(b)  $\sigma_3(z) = \log(1+e^z) \quad \sigma_3'(z) = \frac{e^z}{1+e^z} \quad \forall z \in \mathbb{R} \quad (\text{since } e^z > 0 \quad \forall z)$   
 $= e^z(1+e^z)^{-1}$

$$\sigma_3''(z) = e^z(1+e^z)^{-1} - e^z \cdot e^z(1+e^z)^{-2}$$

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

$$|\sigma_3''(z)| = \left| \frac{e^z}{(1+e^z)^2} \right| = \left| \frac{1}{(1+e^{-z})(1+e^z)} \right| = \left| \frac{1}{2+e^z+e^{-z}} \right| < \frac{1}{2}$$

So  $\forall z \in \mathbb{R}$ ,  $|\sigma_3''(z)|$  is bounded, which implies that  $\sigma_3'(z)$  is Lipschitz continuous.

(for differentiable  $f$ ,  $|f'(x)| \leq M$  for all  $x \in I$  for some  $M > 0 \Rightarrow f$  is L. cts.)

$$\sigma_R(z) = \max\{0, z\} \quad \sigma_R'(z) = \begin{cases} 0 & \text{for } z < 0 \\ 1 & \text{for } z > 0 \\ \text{undefined} & \text{for } z = 0 \end{cases}$$

Assume  $\sigma_R'(z)$  is Lipschitz cts. That is,  $\exists M > 0$  st.  $\forall x, y \in \mathbb{R}$

$$|\sigma_R'(x) - \sigma_R'(y)| \leq M|x-y|$$

Take  $x = \frac{1}{2M}$ ,  $y = -\frac{1}{4M}$ , then

$$|\sigma_R'(x) - \sigma_R'(y)| = |1 - 0| = 1$$

$$\text{and, } M|x-y| = M\left|\frac{1}{2M} + \frac{1}{4M}\right| = \frac{3}{4}$$

$1 \not\leq \frac{3}{4}$  so inequality doesn't hold  $\forall x, y \in \mathbb{R}$  i.e. contradiction  $\Rightarrow \sigma_R'(z)$  is not Lipschitz cts.

$$\sigma(z) = (1 + e^{-z})^{-1}; \rho(z) = (1 - e^{-2z}) / (1 + e^{-2z})$$

5.(c) First note that  $2\sigma(2z) - 1 = \frac{2}{1 + e^{-2z}} - \frac{1 + e^{-2z}}{1 + e^{-2z}} = \frac{1 - e^{-2z}}{1 + e^{-2z}} = \rho(z) \Leftrightarrow \frac{\rho(\frac{1}{2}z) + 1}{2} = \sigma(z)$

We continue by induction.

• For  $L=2$ ,  $y_{2s} = A_2 y_{1s} + b_2$  while  $y_{2t} = C_2 y_{1t} + d_2$   
and  $y_{1s} = \sigma(A_1 x + b_1)$   $y_{1t} = \rho(C_1 x + d_1)$

Using the relationship given above,  $y_{1t} = 2\sigma(2C_1 x + 2d_1) - 1$   $\leftarrow 1 \in \mathbb{R}^n$  of just 1s

$$y_{2t} = 2C_2 \sigma(2C_1 x + 2d_1) - C_2 1 + d_2$$

To represent identical mappings,  $y_{2t} = y_{2s}$ , we would therefore require,

$$2C_2 = A_2; \underline{b}_2 = d_2 - C_2 1; A_1 x + b_1 = 2C_1 x + 2d_1$$

$$\Updownarrow$$

$$A_1 = 2C_1; b_1 = 2d_1$$

So, given  $A_1, A_2, b_1, b_2$ , define,

$$C_1 := \frac{1}{2} A_1 \quad \text{and} \quad d_1 := \frac{1}{2} b_1$$

$$d_2 := b_2 + C_2 1 = b_2 + \frac{1}{2} A_2 1$$

and then the two MLPs represent equivalent mappings with  $y_{2t} = y_{2s}$  for equal  $x$  inputs.

• Now assume true for  $L=i-1$  and consider  $L=i$  (i.e. output of  $(i-1)$  layer networks equivalent),

$$y_{i,s} = A_i y_{i-1,s} + b_i \quad \text{and} \quad y_{i,t} = C_i y_{i-1,t} + d_i$$

By induction step, given  $A_1, \dots, A_{i-1}, b_1, \dots, b_{i-1}$ , it is possible to find  $C_1, \dots, C_{i-1}, d_1, \dots, d_{i-1}$  s.t.

$$y'_{i-1} := A_{i-1} y_{i-2,s} + b_{i-1} = \underbrace{C_{i-1} y_{i-2,t} + d_{i-1}}_{y'_{i-1,t}} \cdot y'_{i-1,t}$$

Hence,  $y_{i-1,s} = \sigma(y'_{i-1})$  and  $y_{i-1,t} = \rho(y'_{i-1})$

$$\Rightarrow y_{i,s} = A_i \sigma(y'_{i-1}) + b_i \quad \text{and} \quad y_{i,t} = C_i \rho(y'_{i-1}) + d_i$$

We desire these to be equal,  $A_i \left( \frac{1}{2} \left( \rho\left(\frac{1}{2} y'_{i-1}\right) + 1 \right) \right) + b_i = C_i \rho(y'_{i-1}) + d_i$

$$\frac{1}{2} A_i \rho\left(\frac{1}{2} y'_{i-1}\right) + \frac{1}{2} A_i 1 + b_i = C_i \rho(y'_{i-1}) + d_i$$

So by setting  $C_i := \frac{1}{2} A_i$ ,  $d_i := \frac{1}{2} A_i 1 + b_i$ , keeping the other inductively given  $C_2, d_2$  but crucially, adjusting  $C_{i-1} \mapsto \frac{1}{2} C_{i-1}$  and  $d_{i-1} \mapsto \frac{1}{2} d_{i-1}$  (so that  $y'_{i-1,t} = \frac{1}{2} y'_{i-1,t}$ ), we can make the final  $i$ -th layer output of the networks too equal (( $i-1$ )-th layers are larger need to be equal).

• Hence, since true for  $L=2$ , true for all  $L \geq 2 \Leftrightarrow$  MLPs are identical mappings. •

(note: if given  $C_1, d_1$  similar rearrangements can be found inductively)

$$\Theta := \{a_1, \dots, a_p, b_1, \dots, b_p, u_1, \dots, u_p\}$$

6. For this optimization problem,  $\underset{\Theta \in \mathbb{R}^{2p}}{\text{minimize}} \frac{1}{N} \sum_{i=1}^N \ell(f_{\Theta}(x_i), y_i)$ ,  
SGD takes the form:

$$\Theta^{k+1} = \Theta^k - r \nabla_{\Theta} \ell_{\Theta}(f_{\Theta}(x_{i(k)}), y_{i(k)}), \quad r \text{ is learning rate}$$

$$\nabla_{\Theta} \ell_{\Theta}(f_{\Theta}(x_{i(k)}), y_{i(k)}) = \nabla_{\Theta} f_{\Theta}(x_{i(k)}) \frac{d\ell_{\Theta}}{dx} + 0$$

$\swarrow$   $y_{i(k)}$  const. wrt  $\Theta$   
 $\nwarrow$  the differential does exist since  $\ell(x, y)$  assumed differentiable in  $x$

$$\nabla_{\Theta} f_{\Theta}(x_{i(k)}) = \nabla_{\Theta} \left( \sum_{j=1}^p u_j \sigma(a_j x_{i(k)} + b_j) \right)$$

$$= \left\{ (\sigma'(a_j x_{i(k)} + b_j) \odot u_j) x_{i(k)}, \sigma'(a_j x_{i(k)} + b_j) \odot u_j, \sigma(a_j x_{i(k)} + b_j) \right\} \text{ by HW 2.6}$$

At initialization,  $k=0$ :

$$\left[ \nabla_{\Theta} f_{\Theta}(x_{i(0)}) \right]_j = \sigma'(a_j^{\circ} x_{i(0)} + b_j^{\circ}) u_j^{\circ} x_{i(0)} = 0 \quad \text{since } a_j^{\circ} x_{i(0)} + b_j^{\circ} < 0 \quad \forall i \text{ (given)}$$

and  $\sigma'(z) = 0 \quad \forall z < 0$

$$\left[ \quad \quad \right]_{p+j} = \sigma'(a_j^{\circ} x_{i(0)} + b_j^{\circ}) u_j^{\circ} = 0$$

$$\left[ \quad \quad \right]_{2p+j} = \sigma(a_j^{\circ} x_{i(0)} + b_j^{\circ}) = 0 \quad \text{by given 'dead condition' directly}$$

So  $j$ -th,  $(p+j)$ -th, and  $(2p+j)$ -th components of  $\nabla_{\Theta} f_{\Theta}(x_{i(0)})$  are 0, so when  $\Theta^{k+1}$  is calculated in SGD these components do not change at all; that is,  $a_j$  and  $b_j$  (and  $u_j$ ) do not change, meaning the above gradient component calculations remain the same (equaling 0) for all  $k \geq 0$  since  $a_j, b_j$  never change from their initialization values which 'killed' the  $j$ -th ReLU layer in the first place. Hence, **if dead on initialization, the  $j$ -th ReLU output remains so throughout learning** i.e. unchanging/invariable

7. If  $\sigma$  is the leaky ReLU function,  $\sigma'(z) = \begin{cases} 1 & \text{for } z \geq 0 \\ \alpha & \text{for } z < 0 \end{cases}$ , i.e.  $\sigma'(z) \neq 0$  for  $z < 0$ .  
So the above become:

$$\left[ \nabla_{\Theta} f_{\Theta}(x_{i(0)}) \right]_j = \sigma'(a_j^{\circ} x_{i(0)} + b_j^{\circ}) u_j^{\circ} x_{i(0)} = \alpha u_j^{\circ} x_{i(0)},$$

$$\left[ \quad \quad \right]_{p+j} = \sigma'(a_j^{\circ} x_{i(0)} + b_j^{\circ}) u_j^{\circ} = \alpha u_j^{\circ},$$

$$\left[ \quad \quad \right]_{2p+j} = \sigma(a_j^{\circ} x_{i(0)} + b_j^{\circ}) = \alpha a_j^{\circ} x_{i(0)} + b_j^{\circ},$$

$\nwarrow$  unexpected factorial :-

none of which are 0! So  $a_j$  and  $b_j$  (and  $u_j$ ) are updated for the next iteration of SGD and the value of **input/output to the  $j$ -th ReLU layer will change** and the gradient always be non-zero (if  $\Theta_i \neq 0 \quad \forall i$ )