

Task 1:

Show that:

$$\max_{\delta \in \Delta} (L(NN_{\theta}(x+\delta), y) - L(NN_{\theta}(x+\delta), y_{\text{target}})) \quad (1)$$

is equal to:

$$\max_{\delta \in \Delta} (NN_{\theta}(x+\delta)_y - NN_{\theta}(x+\delta)_{y_{\text{target}}}) \quad (2)$$

with $NN_{\theta}(x)$ = pre softmax logits

$$\text{and } L(x) = \log\left(\sum_{j=1}^k \exp(NN_{\theta}(x)_j)\right) - NN_{\theta}(x)_y \quad (3)$$

!! stages
a

we just insert (1) into (3) and ignore the max

$$\Rightarrow a - NN_{\theta}(x+\delta)_y - a - NN_{\theta}(x+\delta)_{y_{\text{target}}}$$

we have $a - a$

$$\Rightarrow NN_{\theta}(x+\delta)_y - NN_{\theta}(x+\delta)_{y_{\text{target}}}$$

