

Starting with the loss function:

$$\tilde{J}(\theta) = \underbrace{\frac{1}{N} \sum_{i=1}^N -\log \left[\frac{e^{z_{ij}^{(3)}}}{\sum_k e^{z_{ik}^{(3)}}} \right]}_{J_{\text{data}}} + \underbrace{\lambda (\|w^{(1)}\|_2^2 + \|w^{(4)}\|_2^2)}_{J_{\text{reg}}}$$

Let $\psi(z^{(3)})$ be the probability matrix where each element $\psi(z^{(3)})_{ij}$ is the softmax probability on for class j on sample i :

$$\psi(z^{(3)})_{ij} = \frac{\exp(z_{ij}^{(3)})}{\sum_k \exp(z_{ik}^{(3)})}$$

The loss for each sample i can be expressed as $-\log(\psi(z^{(3)})_{iy_i})$ where y_i is the correct class for sample i :

~~Finding the gradient: (for just~~

$$\frac{\partial \tilde{J}(\theta)}{\partial z_{ij}^{(3)}} = \frac{\partial}{\partial z_{ij}^{(3)}}$$

predicted probability
minus actual
probability

$$\frac{\partial (-\log(\psi(z^{(3)})_{iy_i}))}{\partial (z_{ij}^{(3)})} = \psi(z^{(3)})_{ij} - 1 \quad \text{if } j = y_i$$

and

//

$$= \psi(z^{(3)})_{ij} - 0 \quad \text{if } j \neq y_i$$

$$\text{so } \frac{\partial J_{\text{data}}}{\partial z^{(3)}} = \frac{1}{N} (\psi(z^{(3)}) - \Delta) \quad \text{(putting results in matrix form)}$$

since $\frac{\partial J_{\text{reg}}}{\partial z^{(3)}} = 0$,

$$\frac{\partial J(\theta)}{\partial z^{(3)}} = \frac{1}{N} (\psi(z^{(3)}) - \Delta)$$

Starting with normal loss $J(\theta) = \frac{1}{N} \sum_{i=1}^N -\log \left[\frac{\exp(z_i^{(3)})}{\sum \exp(z_i^{(2)})} \right]$

$$\frac{\partial J(\theta)}{\partial w^{(2)}} = \frac{\partial J}{\partial z^{(3)}} \cdot \frac{\partial z^{(3)}}{\partial w^{(2)}}$$

$$\begin{aligned} \frac{\partial z^{(3)}}{\partial w^{(2)}} &= \frac{\partial}{\partial w^{(2)}} (w^{(2)} a^{(2)} + b^{(2)}) \\ &= a \end{aligned}$$

$$\text{so } \frac{\partial J(\theta)}{\partial w^{(2)}} = \frac{1}{N} (\psi(z^{(3)}) - \Delta) \underbrace{a^{(2)'}}_{\substack{\text{transpose of} \\ a^{(2)} \text{ to make} \\ \text{mult possible}}}$$

Now for $\tilde{J}(\theta) = J(\theta) + J_{\text{reg}}$

$$\begin{aligned} \frac{\partial}{\partial w^{(2)}} (J_{\text{reg}}) &= \frac{\partial}{\partial w^{(2)}} \left(\lambda (\|w^{(1)}\|_2^2 + \|w^{(2)}\|_2^2) \right) \\ &= 2\lambda w^{(2)} \end{aligned}$$

$$\begin{aligned} \text{so } \frac{\partial \tilde{J}(\theta)}{\partial w^{(2)}} &= \frac{\partial J(\theta)}{\partial w^{(2)}} + \frac{\partial J_{\text{reg}}}{\partial w^{(2)}} \\ &= \frac{1}{N} (\psi(z^{(3)}) - \Delta) a^{(2)'} + 2\lambda w^{(2)} \end{aligned}$$

$$\frac{\partial z^{(3)}}{\partial b^{(2)}} = \frac{\partial}{\partial b^{(2)}} (w^{(2)} a^{(2)} + b^{(2)}) = 1$$

$$\begin{aligned} \text{so } \frac{\partial \tilde{J}(\theta)}{\partial b^{(2)}} &= \cancel{\frac{\partial \tilde{J}(\theta)}{\partial z^{(3)}}} \cdot \frac{\partial z^{(3)}}{\partial b^{(2)}} + \frac{\partial J(\theta)}{\partial b^{(2)}} + \frac{\partial J_{\text{reg}}}{\partial b^{(2)}} \\ &= \frac{\partial J(\theta)}{\partial z^{(3)}} \cdot \frac{\partial z^{(3)}}{\partial b^{(2)}} + \frac{\partial}{\partial b^{(2)}} \left(\lambda (\|w^{(1)}\|_2^2 + \|w^{(2)}\|_2^2) \right) \\ &= \frac{1}{N} (\psi(z^{(3)}) - \Delta) (1) + 0 \\ &= \frac{1}{N} (\psi(z^{(3)}) - \Delta) \end{aligned}$$

$$\begin{aligned}
\frac{\partial \tilde{J}(\theta)}{\partial a^{(2)}} &= \frac{\partial J(\theta)}{\partial a^{(2)}} + \frac{\partial J_{reg}}{\partial a^{(2)}} \\
&= \frac{\partial J(\theta)}{\partial z^{(3)}} \cdot \frac{\partial z^{(3)}}{\partial a^{(2)}} + \frac{\partial}{\partial a^{(2)}} (\lambda (\|w^{(1)}\|_2^2 + \|w^{(2)}\|_2^2)) \\
&= \frac{1}{2} (\varphi(z^{(3)}) - \Delta) (w^{(2)})' + 0 \\
&= \frac{1}{2} (\varphi(z^{(3)}) - \Delta) \cdot w^{(2)'} \\
&\quad + 0 \quad \left(\frac{\partial J_{reg}}{\partial a^{(2)}} = 0 \right)
\end{aligned}$$

$$\frac{\partial \tilde{J}(\theta)}{\partial z^{(2)}} = \frac{\partial \tilde{J}(\theta)}{\partial a^{(2)}} \times \frac{\partial a^{(2)}}{\partial z^{(2)}} \quad \checkmark, \text{ since } a^{(2)} = \varphi(z^{(2)})$$

$$\frac{\partial a^{(2)}}{\partial z^{(2)}} = \begin{cases} 1 & \text{if } z^{(2)} > 0 \\ 0 & \text{" } z^{(2)} \leq 0 \end{cases}, \text{ so}$$

$$\frac{\partial \tilde{J}(\theta)}{\partial z^{(2)}} = \begin{cases} \frac{1}{2} (\varphi(z^{(3)}) - \Delta) \cdot w^{(2)'} & \text{if } z^{(2)} > 0 \\ 0 & \text{if } z^{(2)} \leq 0 \end{cases}$$

$$z^{(2)} = w^{(1)} a^{(1)} + b^{(1)}$$

$$\frac{\partial z^{(2)}}{\partial w^{(1)}} = a^{(1)}, \quad \frac{\partial z^{(2)}}{\partial b^{(1)}} = 1$$

$$\frac{\partial J_{reg}}{\partial w^{(1)}} = \frac{\partial}{\partial w^{(1)}} (\lambda (\|w^{(1)}\|_2^2 + \|w^{(2)}\|_2^2)) = 2\lambda w^{(1)}$$

$$\begin{aligned}
\frac{\partial \tilde{J}(\theta)}{\partial w^{(1)}} &= \frac{\partial J(\theta)}{\partial w^{(1)}} + \frac{\partial J_{reg}}{\partial w^{(1)}} \quad (\text{when } z^{(2)} > 0) \\
&= \frac{\partial J(\theta)}{\partial z^{(2)}} \times \frac{\partial z^{(2)}}{\partial w^{(1)}} + 2\lambda w^{(1)} \\
&= \frac{1}{2} (\varphi(z^{(3)}) - \Delta) \cdot a^{(1)'} + 2\lambda w^{(1)}
\end{aligned}$$

Similarly, (when $z^{(2)} > 0$)

$$\begin{aligned}\frac{\partial \tilde{J}(\theta)}{\partial b^{(1)}} &= \frac{\partial J(\theta)}{\partial z^{(2)}} \cdot \frac{\partial z^{(2)}}{\partial b^{(1)}} + \frac{\partial J_{reg}}{\partial b^{(1)}} \\ &= \frac{1}{N} (\varphi(z^{(1)}) - \Delta) \cdot (1) + 0\end{aligned}$$

$$\frac{\partial \tilde{J}(\theta)}{\partial w^{(1)}} = \begin{cases} \frac{1}{N} (\varphi(z^{(1)}) - \Delta) \cdot a^{(1)} & \text{if } z^{(2)} > 0 \\ 0 & \text{if } z^{(2)} \leq 0 \end{cases}$$

$$\frac{\partial \tilde{J}(\theta)}{\partial b^{(1)}} = \begin{cases} \frac{1}{N} (\varphi(z^{(1)}) - \Delta) & \text{if } z^{(2)} > 0 \\ 0 & \text{if } z^{(2)} < 0 \end{cases}$$