---

**Computing derivatives w.r.t Output Layer**

**Part 2**

1. Continuing from where we left off $\dfrac{\partial L(\theta)}{\partial a_{Li}} = \dfrac{-1}{\hat{y}_l} \dfrac{\partial \hat{y}_l}{\partial a_{Li}}$

2. Here, we know that $\hat{y}_l = \dfrac{e^{a_{Ll}}}{\Sigma_i\, a^{Li}}$ (taking the l-th entry of the softmax function applied to vector $a_L$)

3. $\dfrac{\partial L(\theta)}{\partial a_{Li}} = \dfrac{-1}{\hat{y}_l} \dfrac{\partial\, softmax(a_L)_l}{\partial a_{Li}}$ where $a_L = [a_{L1},\, a_{L2} \ldots a_{Lk}]$

   a. Where $softmax(a_L) = [\dfrac{e^{a_{L1}}}{\Sigma_i\, e^{a_{Li}}},\ \dfrac{e^{a_{L2}}}{\Sigma_i\, e^{a_{Li}}} \cdots \dfrac{e^{a_{Lk}}}{\Sigma_i\, e^{a_{Li}}}]$

   b. Selecting the l-th entry would give us the value $softmax(a_L)_l = \dfrac{exp(a_L)_l}{\Sigma_i\, exp(a_L)i}$

4. $\dfrac{\partial L(\theta)}{\partial a_{Li}} = \dfrac{-1}{\hat{y}_l} \dfrac{\partial}{\partial a_{Li}} \dfrac{exp(a_L)_l}{\Sigma_{i'}\, exp(a_L)_{i'}}$

   a. This is of the form $\dfrac{g(x)}{h(x)}$ which gives derivatives $\dfrac{\partial \frac{g(x)}{h(x)}}{\partial x} = \dfrac{\partial g(x)}{\partial x} \dfrac{1}{h(x)} - \dfrac{g(x)}{h(x)^2} \dfrac{\partial h(x)}{\partial x}$

   b. Here $g(x) = exp(a_L)_l$ and $h(x) = \Sigma_{i'}\, exp(a_L)_{i'}$

   c. Substitute the values and expand the formula

5. $\dfrac{\partial L(\theta)}{\partial a_{Li}} = \dfrac{-1}{\hat{y}_l} (\dfrac{\frac{\partial}{\partial a_{Li}} exp(a_L)_l}{\Sigma_{i'}\, exp(a_L)_{i'}} - \dfrac{exp(a_L)_l\, (\frac{\partial}{\partial a_{Li}} \Sigma_{i'}\, exp(a_L)_{i'})}{(\Sigma_{i'}\, exp(a_L)_{i'})^2})$

   a. Here, consider $\dfrac{\partial}{\partial a_{Li}} exp(a_L)_l$ , this value is 0 for all values of i : 0 to k except for when i = l

   b. Thus, we use an indicator variable $1_{(l=i)}\, exp(a_L)_l$ to denote that all other values except i=l resolve to 0

   c. Now consider $\dfrac{\partial}{\partial a_{Li}} \Sigma_{i'}\, exp(a_L)_{i'}$ , here i' ranges from 1 to k. When taking the derivative, only the index i=i' remains, which is simply a derivative of an exponent.

6. $\dfrac{\partial L(\theta)}{\partial a_{Li}} = \dfrac{-1}{\hat{y}_l} (\dfrac{1_{(l=i)}\, exp(a_L)_l}{\Sigma_{i'}\, exp(a_L)_{i'}} - \dfrac{exp(a_L)_l}{\Sigma_{i'}\, exp(a_L)_{i'}} \dfrac{exp(a_L)_i}{\Sigma_{i'}\, exp(a_L)_{i'}})$

   a. This is can be rewritten in terms of the softmax function for the different variables

7. $\dfrac{\partial L(\theta)}{\partial a_{Li}} = \dfrac{-1}{\hat{y}_l} (1_{(l=i)} softmax(a_L)_l - softmax(a_L)_l softmax(a_L)_i$

   a. We know that the Softmax function is $\hat{y}$, so we rewrite it.

8. $\dfrac{\partial L(\theta)}{\partial a_{Li}} = \dfrac{-1}{\hat{y}_l} (1_{(l=i)} \hat{y}_l - \hat{y}_l \hat{y}_i)$

9. After cancellation $\dfrac{\partial L(\theta)}{\partial a_{Li}} = -(1_{(l=i)} - \hat{y}_i)$