Backprop through Softmax based on Input $\frac{\partial L}{\partial a^{[L]}}$

$$\frac{\partial L}{\partial a_b^{[L]}} = -\frac{1}{a_b^{[L]}} \delta_{b,y}$$

$$\frac{\partial L}{\partial z_j^{[L]}} = \sum_b \frac{\partial L}{\partial a_b^{[L]}} \cdot \left( \frac{\partial a_b^{[L]}}{\partial z_j^{[L]}} \right), \quad a_b^{[L]} = \frac{e^{z_b^{[L]}}}{\sum_l e^{z_l^{[L]}}}$$

$$\left[ \frac{\partial a_b^{[L]}}{\partial z_j^{[L]}} = a_b^{[L]} \delta_{bj} - a_b^{[L]} a_j^{[L]} = a_b^{[L]} \left( \delta_{bj} - a_j^{[L]} \right) \right]$$

$$\frac{\partial L}{\partial z_j^{[L]}} = -\sum_b a_b^{[L]} \left( \delta_{bj} - a_j^{[L]} \right) \frac{1}{a_b^{[L]}} \delta_{b,y} = -\left( \delta_{jy} - a_j^{[L]} \right)$$

$$= \left( a_j^{[L]} - \delta_{jy} \right)$$

$$\frac{\partial L}{\partial z_j^{[L]}} = \sum_b \left( a_b^{[L]} \delta_{bj} - a_b^{[L]} a_j^{[L]} \right) \cdot \frac{\partial L}{\partial a_b^{[L]}}$$

$$= a_j^{[L]} \frac{\partial L}{\partial a_j^{[L]}} - \sum_b \left( a_b^{[L]} \frac{\partial L}{\partial a_b^{[L]}} \right) a_j^{[L]}$$

$$\boxed{\frac{\partial L}{\partial z^{[L]}} = a^{[L]} * \frac{\partial L}{\partial a^{[L]}} - \left( a^{[L]} \cdot \frac{\partial L}{\partial a^{[L]}} \right) a^{[L]}}$$

Multiply
element wise

Sum over
the units in
layer, not
over samples
in minibatch

$\frac{\partial L}{\partial a^{[L]}}, \; a^{[L]} \sim (n,m)$

$a^{[L]} \cdot \frac{\partial L}{\partial a^{[L]}} \sim (1,m)$

$\left( sum \left( a^{[L]} * \frac{\partial L}{\partial a^{[L]}}, \; axis=0 \right) \right)$