

$$\frac{\partial C_0}{\partial b^{(l)}} = \frac{\partial z^{(l)}}{\partial b^{(l)}} \cdot \frac{\partial a^{(l)}}{\partial z^{(l)}} \cdot \frac{\partial C_0}{\partial a^{(l)}} = 1 \cdot \text{ReLU}'(z^{(l)}) \cdot 2(a^{(l)} - y)$$

$$4) \frac{\partial z^{(l)}}{\partial b^{(l)}} = 1$$

$$5) \frac{\partial z^{(l)}}{\partial a^{(l)}} = w^{(l)}$$

$$\frac{\partial C_0}{\partial a^{(l-1)}} = \frac{\partial z^{(l)}}{\partial a^{(l-1)}} \cdot \frac{\partial a^{(l)}}{\partial z^{(l)}} \cdot \frac{\partial C_0}{\partial a^{(l)}} = w^{(l)} \cdot \text{ReLU}'(z^{(l)}) \cdot 2(a^{(l)} - y)$$

1) We can only change weights & biases:

$$\frac{\partial C_0}{\partial w^{(l)}} = a^{(l-1)} \cdot \text{ReLU}'(z^{(l)}) \cdot 2(a^{(l)} - y)$$

$$\frac{\partial C_0}{\partial b^{(l)}} = \text{ReLU}'(z^{(l)}) \cdot 2(a^{(l)} - y)$$

So we use the negative of the derivative to stop towards 0. Further we use step of average over a batch of inputs.

2) We can change weights & biases of last layer to change its activation:

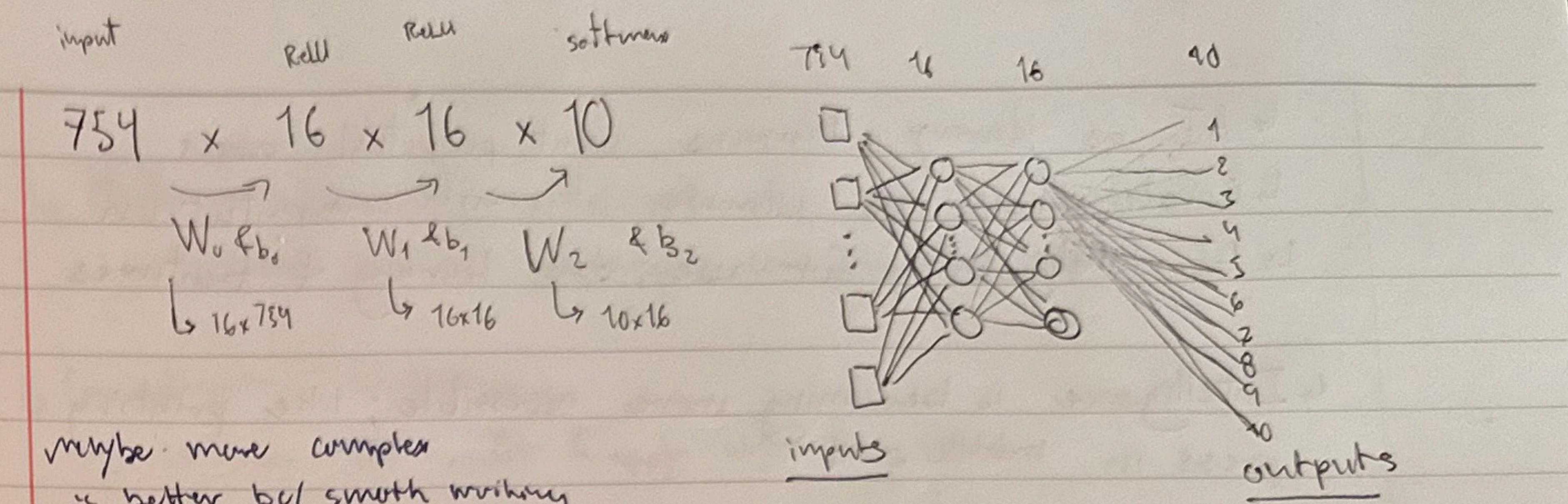
$$\frac{\partial C_0}{\partial w^{(l-1)}} = \frac{\partial z^{(l-1)}}{\partial w^{(l-1)}} \cdot \frac{\partial a^{(l)}}{\partial z^{(l-1)}} \cdot \frac{\partial z^{(l)}}{\partial a^{(l)}} \cdots \frac{\partial C_0}{\partial a^{(l)}} = \frac{\partial a^{(l-1)}}{\partial w^{(l-1)}} \frac{\partial C_0}{\partial a^{(l-1)}}$$

$$a^{(l-1)} = \text{ReLU}(z^{(l-1)}), z^{(l-1)} = w^{(l-1)} a^{(l-2)} + b^{(l-1)}$$

$$\frac{\partial a^{(l-1)}}{\partial w^{(l-1)}} = \frac{\partial z^{(l-1)}}{\partial w^{(l-1)}} \cdot \frac{\partial a^{(l-1)}}{\partial z^{(l-1)}} = a^{(l-2)} \cdot \text{ReLU}'(z^{(l-1)})$$

$$\frac{\partial C_0}{\partial w^{(l-1)}} = a^{(l-2)} \cdot \text{ReLU}'(z^{(l-1)}) \cdot w^{(l)}, \text{ReLU}'(z^{(l)}) \cdot 2(a^{(l)} - y)$$

stored from earlier



$\hat{a}^i \rightarrow$  determined by image inputs

$a^0 = \text{ReLU}(W_0 \hat{a}^0 + b_0) \rightarrow$  \* we could batch norm? (no later)

$a^1 = \text{ReLU}(W_1 a^0 + b_1)$

$a^2 = \text{softmax}(W_2 a^1 + b_2)$  (where  $y$  is desired output)

$$\text{Loss } C_0 = \frac{1}{N} [(y - \hat{a}^2)^2] [1] = 1 \times 1 \text{ cost matrix}$$

$$\text{ReLU}(x) = \begin{cases} 0 & \text{if } x \leq 0 \\ x & \text{if } x > 0 \end{cases}, \text{ReLU}' = \begin{cases} 0 & \text{if } x \leq 0 \\ 1 & \text{if } x > 0 \end{cases}$$

$$z^{(l)} = w^{(l)} a^{(l-1)} + b^{(l)}, a^{(l)} = \text{ReLU}(z^{(l)}), \hat{a}^2 = (a^{(l)} - y)^2$$

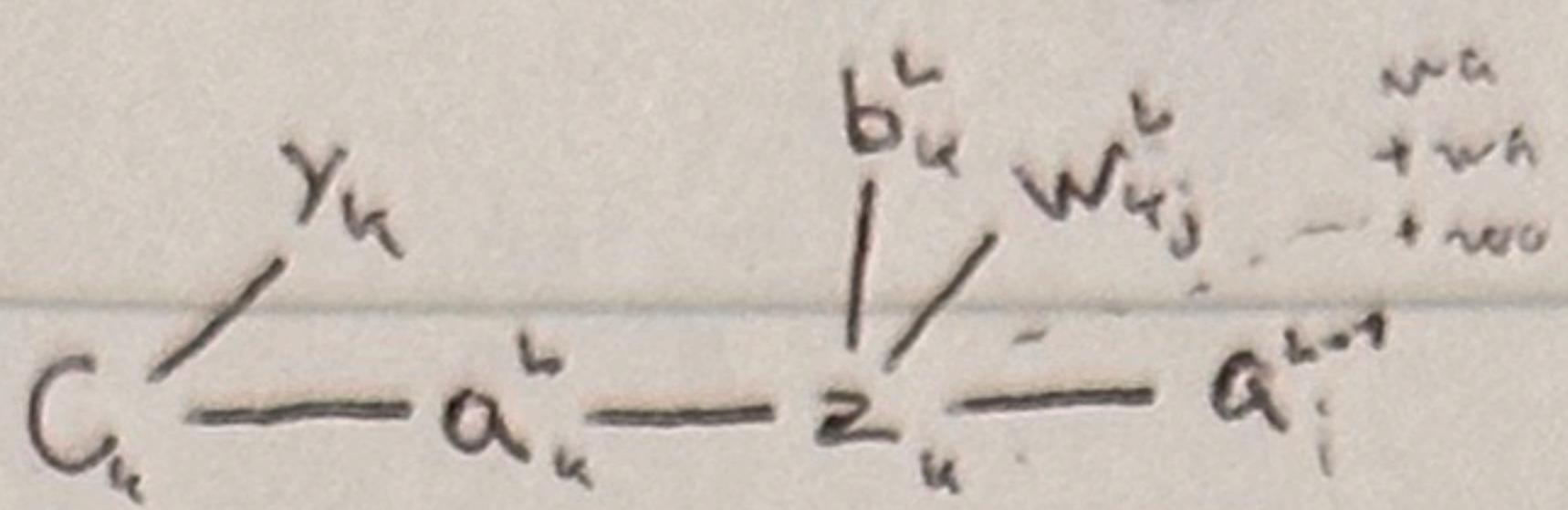
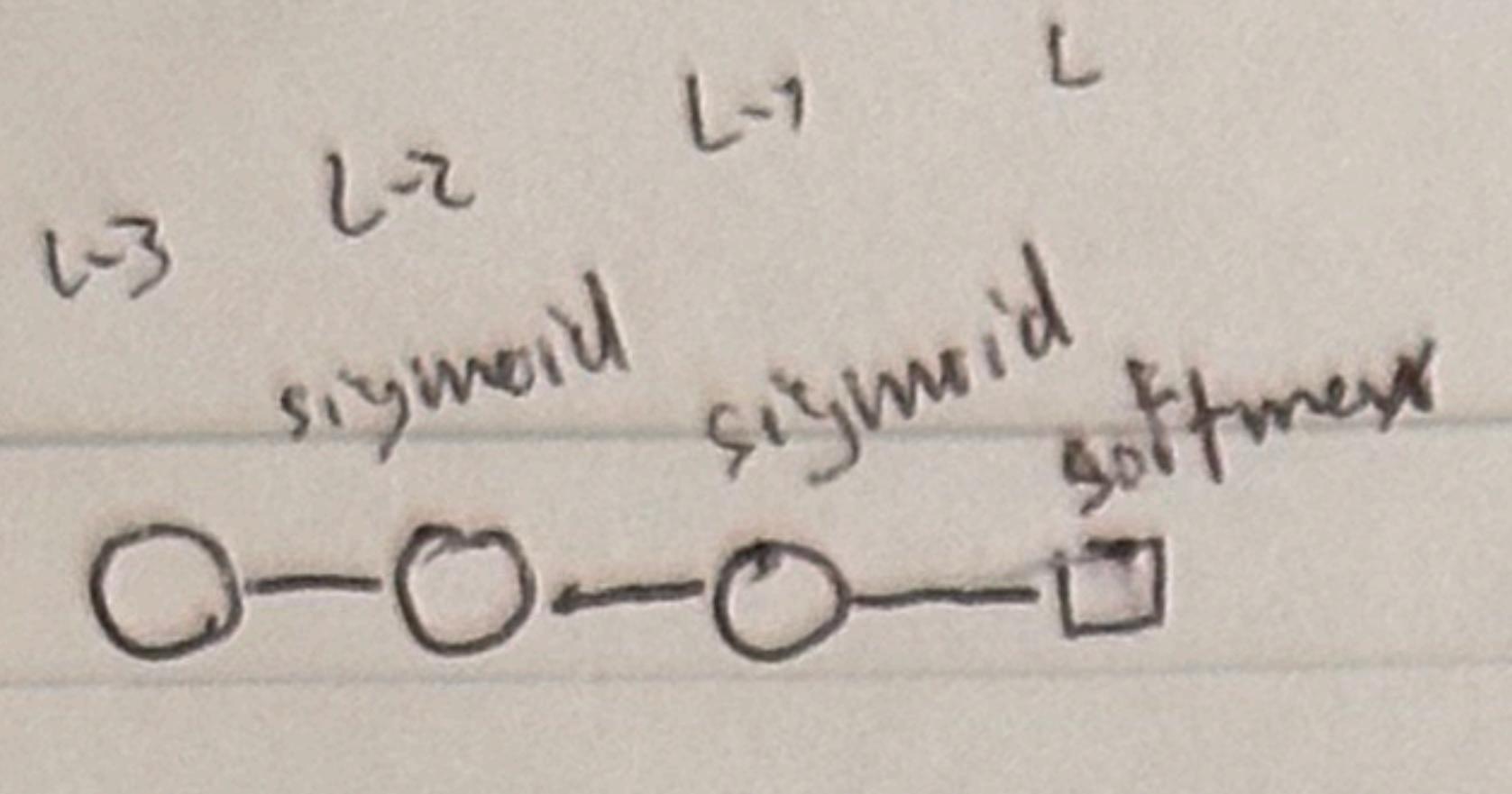
We want sensitivity of  $C_0$  to change in 1)  $w^{(l)}$  2)  $b^{(l)}$  3)  $a^{(l-1)}$

$$\frac{\partial C_0}{\partial w^{(l)}} = \frac{\partial z^{(l)}}{\partial w^{(l)}} \frac{\partial a^{(l)}}{\partial z^{(l)}} \frac{\partial C_0}{\partial a^{(l)}}$$

$$1) \frac{\partial C_0}{\partial a^{(l)}} = 2(a^{(l)} - y) \quad 2) \frac{\partial a^{(l)}}{\partial z^{(l)}} = \text{ReLU}'(z^{(l)})$$

$$3) \frac{\partial z^{(l)}}{\partial w^{(l)}} = \cancel{w^{(l)}} a^{(l-1)}$$

$$\therefore \text{Thus } \frac{\partial C_0}{\partial w^{(l)}} = a^{(l-1)} \text{ReLU}'(z^{(l)}) 2(a^{(l)} - y)$$



C is cross entropy loss

$$\text{sigmoid} = \sigma(z_k^l) = \frac{1}{1 + e^{-(z_k^l)}}$$

$$\Rightarrow a_k^l = \sigma(z_k^l)$$

$$\text{softmax} = \sigma_z(z_k^l) = \frac{e^{z_k^l}}{\sum_{k=1}^K e^{z_k^l}}$$

$$\text{sigmoid}' = \sigma(z_k^l)(1 - \sigma(z_k^l)) = a_k^l(1 - a_k^l)$$

$$\text{softmax}' = \sigma_z'(z_k^l) \quad ?$$

todo

~~(1)  $\frac{\partial C}{\partial a_k^l} = \sum_{j=1}^K \frac{\partial C}{\partial a_j^l} \frac{\partial a_j^l}{\partial a_k^l}$~~

$$\frac{\partial C}{\partial a_k^l} = -\frac{t_k}{a_k^l}$$

where  $t_k$  is target  
and  $a_k$  is hypothesis

~~(2)  $a_k^l = \frac{e^{z_k^l}}{\sum_{k=1}^K e^{z_k^l}}$~~

~~(4)  $C = -\sum_{k=1}^K t_k \log(a_k^l)$~~

~~(3)  $\frac{\partial C}{\partial w_{kj}^l} = \frac{\partial C}{\partial a_k^l} = \left( \sum_{\text{output neurons}} \frac{\partial C}{\partial a_k^l} \frac{\partial a_k^l}{\partial z_j^l} \right) \frac{\partial z_j^l}{\partial w_{kj}^l}$~~

$$\frac{\partial a_k^l}{\partial z_j^l} = \begin{cases} a_j^l(1 - a_j^l) & \text{if } k=j \\ -a_j^l a_k^l & \text{if } k \neq j \end{cases}$$

~~(5)  $\sum_{\text{neurons}} \frac{\partial C}{\partial a_k^l} \frac{\partial a_k^l}{\partial z_j^l} = \sum_{\text{neurons}, k \neq j} a_j^l t_j = a_k^l - t_k$~~

$$\frac{\partial C}{\partial z_j^l}$$

Thus (5)  $\frac{\partial C}{\partial b_k^l} = \frac{\partial C}{\partial z_k^l} \frac{\partial z_k^l}{\partial b_k^l} = (a_k^l - t_k) \cdot 1 = (a_k^l - t_k)$

$$\frac{\partial z_k^l}{\partial a_k^l} \frac{\partial C}{\partial z_k^l}$$

~~(6)  $\frac{\partial C}{\partial w_{kj}^l} = \frac{\partial z_k^l}{\partial w_{kj}^l} \frac{\partial C}{\partial z_k^l} = a_j^{(l-1)} \cdot (a_k^l - t_k)$~~

~~(7)  $\frac{\partial C}{\partial a_j^{(l-1)}} = \sum_{k=1}^K w_{kj}^{(l)} \cdot (a_k^l - t_k)$~~

~~(1)  $a_j^{(l-1)} = \sigma(z_j^{(l-1)}) = \frac{1}{1 + e^{-(z_j^{(l-1)})}}$~~

\* On input layer use input values.

\* Not done on front most layer

~~(2)  $\frac{\partial C}{\partial b_j^{(l-1)}} = \frac{\partial z_j^{(l-1)}}{\partial b_j^{(l-1)}} \frac{\partial a_j^{(l-1)}}{\partial z_j^{(l-1)}} \frac{\partial C}{\partial a_j^{(l-1)}} = 1 \cdot a_j^{(l-1)}(1 - a_j^{(l-1)}) \frac{\partial C}{\partial a_j^{(l-1)}}$~~

~~(3)  $\frac{\partial C}{\partial w_{jh}^{(l-1)}} = \frac{\partial z_j^{(l-1)}}{\partial w_{jh}^{(l-1)}} \frac{\partial a_j^{(l-1)}}{\partial z_j^{(l-1)}} \frac{\partial C}{\partial a_j^{(l-1)}} = a_h^{(l-2)} \cdot a_j^{(l-1)}(1 - a_j^{(l-1)}) \frac{\partial C}{\partial a_j^{(l-1)}}$~~

~~(4)  $\frac{\partial C}{\partial a_h^{(l-2)}} = \sum_{j=1}^J \frac{\partial z_j^{(l-1)}}{\partial a_h^{(l-2)}} \frac{\partial a_j^{(l-1)}}{\partial z_j^{(l-1)}} \frac{\partial C}{\partial a_j^{(l-1)}} = \sum_{j=1}^J w_{jh}^{(l-1)} \cdot a_j^{(l-1)}(1 - a_j^{(l-1)}) \frac{\partial C}{\partial a_j^{(l-1)}}$~~

Sigmoid

$$e(\Sigma - c)$$

$$C_o = \sum_{k=1}^K (a_k^{(l)} - y_k)^2, a_k^{(l)} = \text{softmax}(z_k^{(l)}) = \frac{e^{z_k^{(l)}}}{\sum_{k=1}^K e^{z_k^{(l)}}}$$

$$\therefore C_o = \left[ \sum_{k=1}^K \frac{e^{z_k^{(l)}}}{\sum_{k=1}^K e^{z_k^{(l)}}} - y_k \right]^2 \quad \left\{ C_o \text{ in terms of } z_k^{(l)} \right.$$

$$\frac{\partial C_o}{\partial w_{kj}^{(l)}} = \frac{\partial z_k^{(l)}}{\partial w_{kj}^{(l)}} \frac{\partial a_k^{(l)}}{\partial z_k^{(l)}} \frac{\partial C_o}{\partial a_k^{(l)}} = a_j^{(l-1)} \cdot \text{softmax}(z_k^{(l)}) \cdot (1 - \text{softmax}(z_k^{(l)})) \cdot 2(a_k^{(l)} - y_k)$$

$$\text{or} = \frac{\partial z_k^{(l)}}{\partial w_{kj}^{(l)}} \cdot \frac{\partial C_o}{\partial z_k^{(l)}} = a_j^{(l-1)} \cdot \left[ \underbrace{g(z_k^{(l)}) (1 - g(z_k^{(l)})) \cdot 2(a_k^{(l)} - y_k) + \dots + -\sigma(z_h^{(l)}) \sigma(z_k^{(l)}) 2(a_k^{(l)} - y_k)}_{\text{activation}} \right] \underbrace{-\sigma(z_h^{(l)}) \sigma(z_k^{(l)}) 2(a_k^{(l)} - y_k)}_{\text{derivative}}$$

\* right

- 1) Forward pass & store activations in vectors for each layer } after activation function  
 2) calculate derivative of loss w.r.t. expected each neuron & store in vector  $\langle 2(a_k^{(l)} - y_k) \rangle$  } step 0(l)  
 3.1) Mult by neuron's activation (this is  $a_k^{(l)} \cdot 2(a_k^{(l)} - y_k)$ ) store in vector } step 1(l)  
 3.2) Now for each neuron mult its activation by  $(1 - a_k^{(l)})$  add  $a_k^{(l)} \cdot 2(a_k^{(l)} - y_k)$  adding of all other vectors } step 2(l)  
 3) do  $\frac{\partial C_o}{\partial w_{kj}^{(l)}} = a_j^{(l-1)} \cdot (\text{activation}_k)$  (store gradient +=)  
 4) do  $\frac{\partial C_o}{\partial b_k^{(l)}} = \frac{\partial z_k^{(l)}}{\partial b_k^{(l)}} \cdot \frac{\partial C_o}{\partial z_k^{(l)}} = 1 \cdot (\text{activation}_k)$  store bias gradient += } step 0(l)  
 5) Contribute to activation bias of next layer  $\frac{\partial C_o}{\partial a_j^{(l+1)}} = \frac{\partial z_k^{(l)}}{\partial a_j^{(l+1)}} \frac{\partial C_o}{\partial z_k^{(l)}} = w_{kj}^{(l)} \cdot (\text{activation}_k)$  } store in same array NOT  
 6) same thing 1) bias 2) weights 3) contribute to next  
 7) at least use no #3, make sure to use clear/formatted arrays

1) negative activation of self    2) add  $(1 - a_k^{(l)})$  div 1

repurpose to  $\frac{\partial C_o}{\partial a_k^{(l)}}$   
activation 1

### Activation Functions:

$$\text{ReLU} = \begin{cases} 0 & \text{if } z \le 0 \\ z & \text{if } z > 0 \end{cases}, \quad \text{ReLU}' = \begin{cases} 0 & \text{if } z \le 0 \\ 1 & \text{if } z > 0 \end{cases}$$

$$\text{softmax} = \frac{e^{z_k^{(l)}}}{h^{(l)}}, \quad \text{softmax}(z_k^{(l)}) = a_k^{(l)} = \frac{e^{z_k^{(l)}}}{\sum_{k=1}^K e^{z_k^{(l)}}}$$

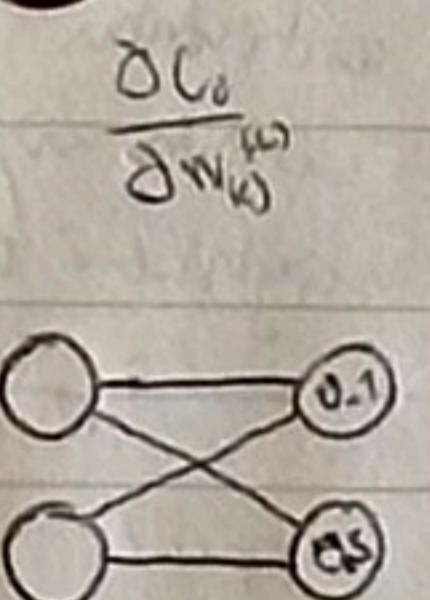
$$h^{(l)} = \sum_{k=0}^{n-1} e^{z_k^{(l)}} = e^{z_0^{(l)}} + e^{z_1^{(l)}} + e^{z_2^{(l)}} + \dots + e^{z_K^{(l)}}$$

$$\frac{d}{dx} \frac{f(x)}{g(x)} = \frac{f'(x)g(x) - f(x)g'(x)}{(g(x))^2}$$

$$= \frac{e^{z_k^{(l)}}}{\sum_{k=1}^K e^{z_k^{(l)}}} - \frac{e^{z_k^{(l)}}}{\sum_{k=1}^K e^{z_k^{(l)}}} e^{z_k^{(l)}}$$

$$a_k^{(l)} = \text{softmax}(z_k^{(l)}), \quad \frac{\partial a_k^{(l)}}{\partial z_k^{(l)}} = \frac{\partial}{\partial z_k^{(l)}} \frac{e^{z_k^{(l)}}}{\sum_{k=1}^K e^{z_k^{(l)}}} = \frac{e^{z_k^{(l)}}}{\sum_{k=1}^K e^{z_k^{(l)}}} \left[ \sum_{k=1}^K e^{z_k^{(l)}} \right]^2$$

$$\frac{\partial C_o}{\partial w_{kj}^{(l)}} = \frac{\partial z_k^{(l)}}{\partial w_{kj}^{(l)}} \frac{\partial a_k^{(l)}}{\partial z_k^{(l)}} \frac{\partial C_o}{\partial a_k^{(l)}} = \frac{e^{z_k^{(l)}}}{\sum_{k=1}^K e^{z_k^{(l)}}} \left( 1 - \frac{e^{z_k^{(l)}}}{\sum_{k=1}^K e^{z_k^{(l)}}} \right) \cdot a_j^{(l-1)} \cdot 2(a_k^{(l)} - y_k) = \frac{e^{z_k^{(l)}}}{\sum_{k=1}^K e^{z_k^{(l)}}} \left( 1 - \frac{e^{z_k^{(l)}}}{\sum_{k=1}^K e^{z_k^{(l)}}} \right) \text{sum} = a_k^{(l)} (1 - \sigma(z_k^{(l)}))$$



↳ I need to consider overall cost of changing weights;  
this means I need to sum & check others

where:  $h \neq k$

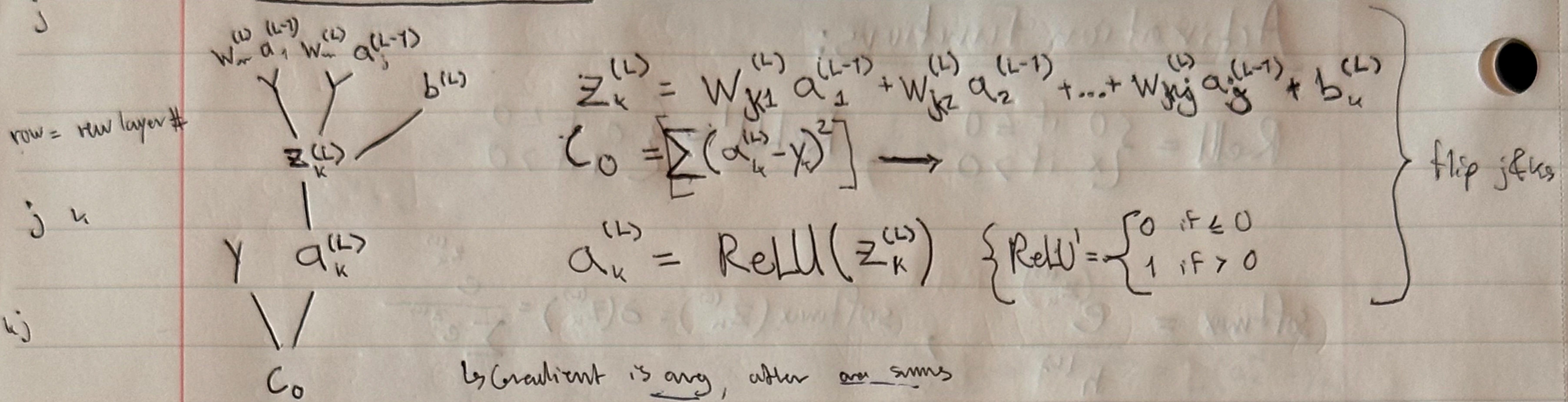
$$\frac{\partial C_o}{\partial a_j^{(l+1)}} = \frac{\partial}{\partial a_j^{(l+1)}} \left[ \frac{e^{z_h^{(l)}}}{\sum_{k=1}^K e^{z_k^{(l)}}} \right] = \frac{\partial}{\partial z_h^{(l)}} \left[ \frac{e^{z_h^{(l)}}}{\sum_{k=1}^K e^{z_k^{(l)}}} \right]$$

$$\frac{e^{z_h^{(l)}}}{e^{z_h^{(l)}} + e^{z_k^{(l)}}} = \frac{0 \cdot \sum_{k=1}^K e^{z_k^{(l)}} - e^{z_h^{(l)}} e^{z_k^{(l)}}}{\left( \sum_{k=1}^K e^{z_k^{(l)}} \right)^2} = -\sigma(e^{z_h^{(l)}}) \sigma(e^{z_k^{(l)}})$$

$W \in K \times J$  matrix

$W_{jk}^{(l)}$   $\hookrightarrow$  connection, edge in Layer L-1  
 $\hookrightarrow$  current node in layer L

### Additional Neurons:



many per node

$\frac{\partial C_0}{\partial W_{kj}^{(L)}} = \frac{\partial Z_k^{(L)}}{\partial W_{kj}^{(L)}} \frac{\partial a_j^{(L-1)}}{\partial Z_k^{(L)}} \frac{\partial C_0}{\partial a_j^{(L-1)}} = a_j^{(L-1)} \cdot \text{ReLU}'(Z_k^{(L)}) \cdot 2(a_k^{(L)} - y_k)$

many per layer

only once

$\frac{\partial C_0}{\partial b_k^{(L)}} = \frac{\partial Z_k^{(L)}}{\partial b_k^{(L)}} \frac{\partial a_k^{(L-1)}}{\partial Z_k^{(L)}} \frac{\partial C_0}{\partial a_k^{(L-1)}} = 1 \cdot \text{ReLU}'(Z_k^{(L)}) \cdot 2(a_k^{(L)} - y_k)$

many per node

many per layer

only for  $k=3$  node

for real activation sum all weights w/ other neighbor nodes

$\frac{\partial C_0}{\partial a_j^{(L-1)}} = \frac{\partial Z_k^{(L)}}{\partial a_j^{(L-1)}} \frac{\partial a_k^{(L-1)}}{\partial Z_k^{(L)}} \frac{\partial C_0}{\partial a_k^{(L-1)}} = W_{kj}^{(L)} \cdot \text{ReLU}'(Z_k^{(L)}) \cdot 2(a_k^{(L)} - y_k)$

abc repeat for all edges

$\frac{\partial C_0}{\partial b_j^{(L-1)}} = \frac{\partial Z_k^{(L-1)}}{\partial b_j^{(L-1)}} \frac{\partial a_k^{(L-1)}}{\partial Z_k^{(L-1)}} \frac{\partial C_0}{\partial a_k^{(L-1)}} = 1 \cdot \text{ReLU}'(Z_k^{(L-1)}) \cdot \boxed{\frac{\partial C_0}{\partial a_j^{(L-1)}}}$

$\frac{\partial C_0}{\partial a_i^{(L-2)}} = \frac{\partial Z_k^{(L-1)}}{\partial a_i^{(L-2)}} \frac{\partial a_j^{(L-1)}}{\partial Z_k^{(L-1)}} \frac{\partial C_0}{\partial a_j^{(L-1)}} = W_{ji}^{(L-1)} \cdot \text{ReLU}'(Z_k^{(L-1)}) \cdot \boxed{\frac{\partial C_0}{\partial a_j^{(L-1)}}}$

↳ For last layer (frontmost) - don't calc prev layer's influence  
 ↳ At end we have gradient for biases & weights; these accumulate & divide at end by batch size