

## Question

①

How much does  $J(X, Y)$  change along a dimension  $\theta^{[l]}$  in the space  $(W, B)$ , the space of all weights and bias values, at point  $K$ , a point in the larger space  $\leftarrow$  larger by  $m \cdot (\lambda^{[0]} + 1)$  dimensions.  $(W, B, X, Y)$  of weights, biases, training example inputs, training example label.

$J(X, Y)$  has been computed at point  $K \in (W, B, X, Y)$

## Find Gradient

We need to find the derivatives

$$\left. \frac{\partial J}{\partial \theta^{[l]}} \right|_K$$

Notation for: the partial derivative of  $J$  along the dimension of "placeholder variable"  $\theta^{[l]}$ , which is related to layer  $[l]$  and will be further defined, at point  $K$ .

$\theta^{[l]}$  is one of the weights or bias values found in layer  $[l]$ .

$$\theta^{[l]}$$

is one of:

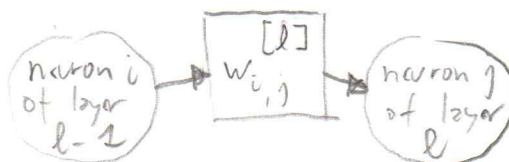
$$w_{i,j}^{[l]}$$

$l \in \{1, 2, 3\}$  for a 3-layer net

$$i \in \{1, \dots, 2^{[l-1]}\}$$

$$j \in \{1, \dots, 2^{[l]}\}$$

$$b_j^{[l]}$$



$$j \in \{1, \dots, 2^{[l]}\}$$

(2)

$$\left. \frac{\partial J}{\partial \theta^{[L]}} \right|_K = \frac{1}{m} \sum_{e \in \{1, \dots, m\}} \left( \frac{\partial}{\partial \theta^{[L]}} \text{loss}(\hat{y}^{(e)}, y^{(e)}) \right) \Big|_K$$

this  
computes  
a mean

In the diagram, the above values are arranged  
in a matrix according to the  $\theta^{[L]}$ .

Developing only the expression of the term :

- 1) use our definition of  $\text{loss}(-, -)$   
based on negative log-likelihood for 2 classes
- 2)  $y^{(e)}$  does not depend on  $\theta^{[L]}$

$$\begin{aligned} & - \left[ y^{(e)} \frac{\partial}{\partial \theta^{[L]}} \log(\hat{y}^{(e)}) + (1 - y^{(e)}) \frac{\partial}{\partial \theta^{[L]}} \log(1 - \hat{y}^{(e)}) \right] \\ & \quad \quad \quad \swarrow \text{chain rule} \quad \quad \quad \searrow \text{chain rule} \\ & = - \left[ y^{(e)} \cdot \frac{1}{\hat{y}^{(e)}} \cdot \frac{\partial}{\partial \theta^{[L]}} \hat{y}^{(e)} + (1 - y^{(e)}) \frac{1}{1 - \hat{y}^{(e)}} \cdot \frac{\partial}{\partial \theta^{[L]}} (1 - \hat{y}^{(e)}) \right] \end{aligned}$$

Using  $\hat{y}^{(e)} = a^{[3]}(e)$  ( $a^{[3]}(e)$  is of the shape  $1 \times 1$  i.e. a scalar!)

Pulling in the minus

Differentiating the rightmost sum

$$= - \frac{y^{(e)}}{\hat{y}^{(e)}} \cdot \frac{\partial}{\partial \theta^{[L]}} a^{[3]}(e) - \frac{1 - y^{(e)}}{1 - \hat{y}^{(e)}} \cdot \frac{\partial}{\partial \theta^{[L]}} (1 - a^{[3]}(e))$$

$$= \left( \frac{1 - y^{(e)}}{1 - \hat{y}^{(e)}} - \frac{y^{(e)}}{\hat{y}^{(e)}} \right) \cdot \frac{\partial}{\partial \theta^{[L]}} a^{[3]}(e)$$

scalar in this architecture!

$$\left. \frac{\partial}{\partial \eta} \text{loss}(\eta, y^{(e)}) \right|_{\eta = \hat{y}^{(e)}}$$

Note that in the expression

$$\frac{1 - y^{(c)}}{1 - \hat{y}^{(c)}} - \frac{y^{(c)}}{\hat{y}^{(c)}}$$

$y^{(c)}$  and  $1 - y^{(c)}$  are "selectors": exactly one is 1.

This can be used during implementations to drop one of the terms and avoid floating-point division that may yield bad results. One should probably make sure to clamp the division results to reasonable values, too.

Now we need to work on the "red block" from earlier:

$$\frac{\partial}{\partial \phi^{[l]}} a^{[3]}(e)$$

For layer [3], we use the standard sigmoid  $\sigma(-)$  as activation function!

$$= \frac{\partial}{\partial \phi^{[l]}} \sigma(z^{[3]}(e))$$

Note that for layer [3]  
 $z^{[3]}(e)$  has shape  $1 \times 1$  (a scalar)

$$= \underbrace{\frac{\partial}{\partial \eta} \sigma(\eta)}_{\left|_{z^{[3]}(e)}\right.} \cdot \underbrace{\frac{\partial}{\partial \phi^{[l]}} z^{[3]}(e)}$$

(chain rule, introduce a new variable  $\eta$  for the expression of the derivative of  $z^{[3]}(e)$ )

$$= \sigma' \left( z^{[3]}(e) \right) \cdot \frac{\partial}{\partial g^{[3]}} z^{[3]}(e)$$

Switch to Lagrange's  
"prime notation"

Now use a property peculiar to  $\sigma$ :  $\sigma' = \sigma \circ (1 - \sigma)$

$$= \underbrace{\sigma \left( z^{[3]}(e) \right) \cdot \left( 1 - \sigma \left( z^{[3]}(e) \right) \right)}_{\text{this part can be computed and cached during forward computation}} \cdot \frac{\partial}{\partial g^{[3]}} z^{[3]}(e)$$

this part can be  
computed and cached during  
forward computation


Now we just need to work on  
above.

the "red block"

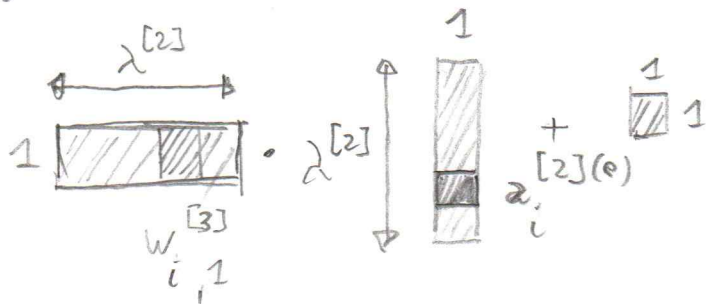
How does it look for layer 3?

$$\frac{\partial}{\partial g^{[3]}} z^{[3]}(e)$$

NB now at layer 3

$z^{[3]}(e)$  is  $1 \times 1$  

$$= \frac{\partial}{\partial g^{[3]}} \left( W^{[3]T} \cdot a^{[2]}(e) + b^{[3]} \right)$$



► For  $g^{[3]} = w_{i,j}^{[3]}$

neuron of  
layer 2

neuron of  
layer 3, can  
only be 1

$$\frac{\partial}{\partial w_{i,1}^{[3]}} \left( W^{[3]T} \cdot a^{[2]}(e) + b^{[3]} \right) = a_i^{[2]}(e)$$

(5)

For  $p^{[3]} = b_1^{[3]}$  :  $\frac{\partial}{\partial b_1^{[3]}} \left( \underbrace{W^{[3]T} \cdot a^{[2]}(e)}_{\rightarrow 0} + \underbrace{b_1^{[3]}}_{\rightarrow 1} \right)$

again, the index of the neuron for layer [3] can only be 1 ...

=  $\boxed{1}$

This ends the calculation for layer [3] as we now have all the rfs:

$$\left. \frac{\partial J}{\partial p^{[3]}} \right|_{\mathbb{R}} = \frac{1}{m} \sum_{e \in \{1, \dots, m\}} \left[ \left( \frac{1 - y^{(e)}}{1 - \hat{y}^{(e)}} - \frac{y^{(e)}}{\hat{y}^{(e)}} \right) \right. \quad \text{factor from loss}$$

$$\cdot \sigma(z_1^{[3]}(e)) \cdot (1 - \sigma(z_1^{[3]}(e))) \quad \text{factor from act}^{[3]}$$

$$\cdot \left. \begin{cases} p^{[3]} = b_1^{[3]} & ; & 1 \\ p^{[3]} = w_{i,1}^{[3]} & ; & a_i^{[2]}(e) \end{cases} \right]$$



6

How does it look for layer 2?

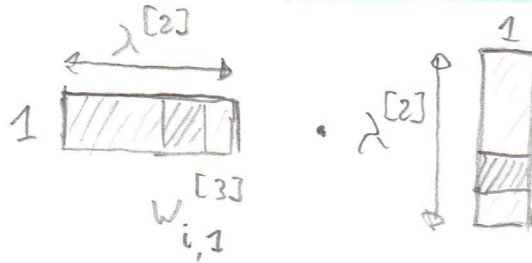
$$\frac{\delta}{\delta p^{[2]}} z^{[3]}(e)$$

NB now at layer 2

$$= \frac{\delta}{\delta p^{[2]}} \left( W^{[3]T} \cdot a^{[2]}(e) + b^{[3]} \right)$$

resolves to 0 when differentiating

$$= W^{[3]T} \cdot \frac{\delta}{\delta p^{[2]}} a^{[2]}(e)$$



Now we need to work on the "green block" above:

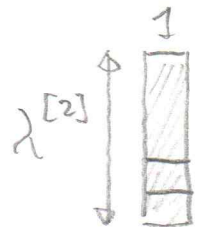
$$\frac{\delta}{\delta p^{[2]}} a^{[2]}(e)$$

For layer [2], we use an unspecified activation function  $act^{[2]}(-)$ ;

$$= \frac{\partial}{\partial p^{[2]}} act^{[2]}(z^{[2]}(e))$$

$$= \frac{\partial}{\partial \eta} act^{[2]}(\eta) \bigg|_{z^{[2]}(e)} \cdot \frac{\partial}{\partial p^{[2]}} z^{[2]}(e)$$

element wise multiplication



Shape of  $z^{[2]}(e)$ .  
The derivative has the same shape

New variable  $\eta$  for clarity

⑦

$$= \underbrace{\text{act}^{[2]'}(z^{[2]}(e))}_{\text{can be computed and cached during fwd computation}} \cdot \frac{\partial}{\partial \phi^{[2]}} z^{[2]}(e)$$

Switch to  
Lagrange's  
"prime notation"

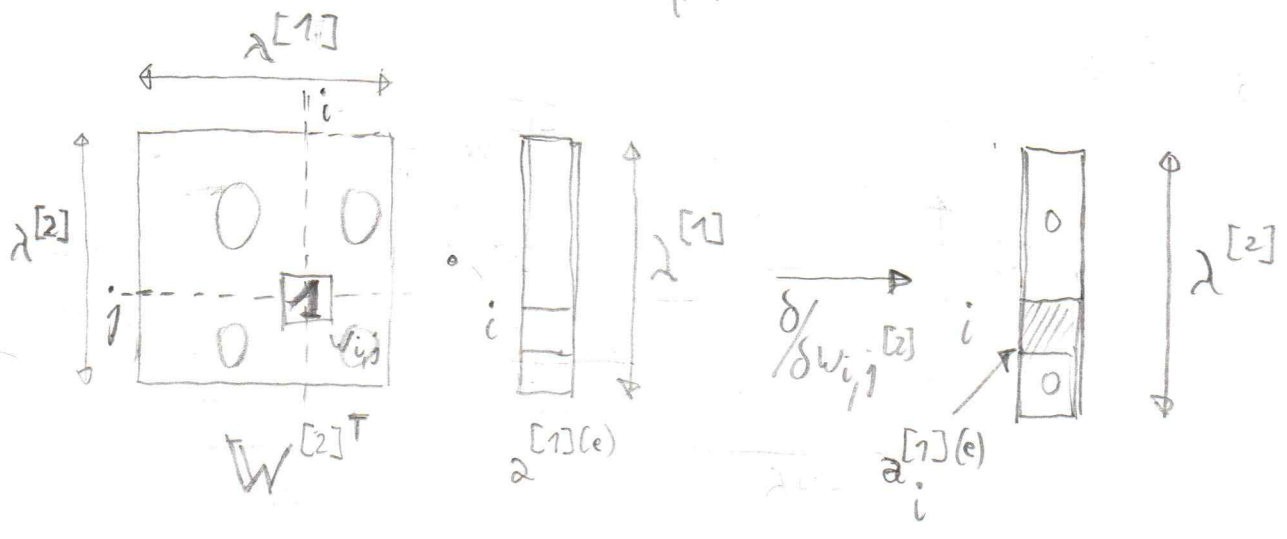
Now we just need to work on the "red block" above:

$$\frac{\partial}{\partial \phi^{[2]}} z^{[2]}(e) = \frac{\partial}{\partial \phi^{[2]}} \left( W^{[2]T} \cdot a^{[1]}(e) + b^{[2]} \right)$$

► For  $\phi^{[2]} = w_{i,j}^{[2]}$ :

$$\frac{\partial}{\partial w_{i,j}^{[2]}} \left( W^{[2]T} \cdot a^{[1]}(e) + b^{[2]} \right) = a_i^{[1]}(e)$$

From neuron  $i$  in  $[1]$  to neuron  $j$  in  $[2]$  selects  $a_i^{[1]}(e)$  on differentiation yields  $\delta$  on differentiation inside a column vector of  $\phi$ s



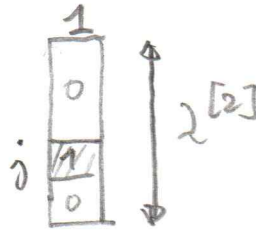
(8)

For  $\delta^{[2]} = b_j^{[2]}$  :

$$\frac{\delta}{\delta L_j^{[2]}} \left( \underbrace{W^{[2]T} \cdot a^{[1]}(e)}_{\text{yields } 0} + \underbrace{b^{[2]}}_{\text{yields "1 at place j"}} \right) = 1$$

inside a  
column  
vector of  $\delta s$

the above has the shape



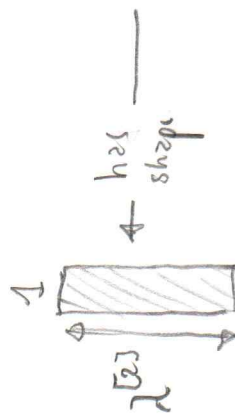
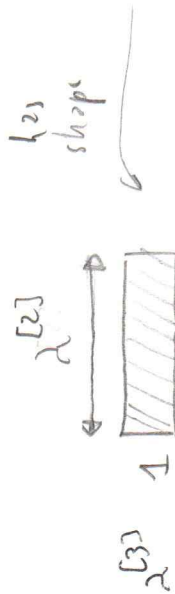
This ends the calculation for layer [2] as we now have all the info.



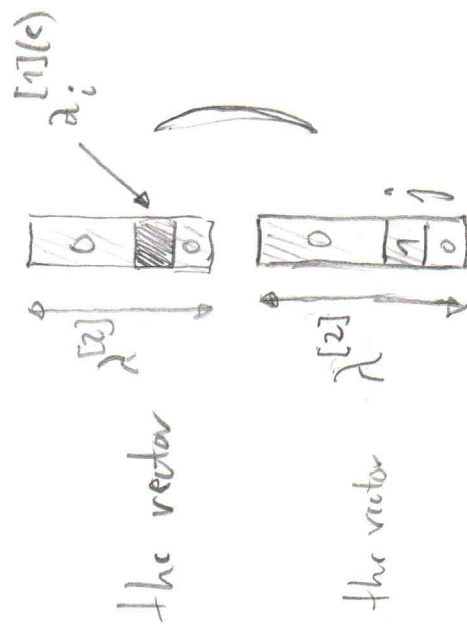
9

$$\frac{\partial J}{\partial g^{[2]}} \Big|_K = \frac{1}{m} \sum_{e \in \{1, \dots, m\}}$$

For a given  $g^{[2]}$ ,  
this is a scalar.



$$\left\{ \begin{array}{l} \text{if } g^{[2]} = w_{i,j}^{[2]} \\ \text{if } g^{[2]} = b_j^{[2]} \end{array} \right.$$



the vector

the vector

9

scalar factor  
from loss computation

$$\left( \frac{1-y^{(e)}}{1-y^{(e)}} - \frac{y^{(e)}}{y^{(e)}} \right)$$

scalar factor  
from differentiating  
through  $\text{act}^{[3]} = \sigma(-)$

$$\sigma'(z^{[3]}(e)) \cdot (1 - \sigma(z^{[3]}(e)))$$

from differentiating  
 $z^{[3]}(e)$

$$W^{[3]T}$$

vector factor  
from differentiating  
through  $\text{act}^{[2]}$

$$\text{act}^{[2]'}(z^{[2]}(e))$$

And J' ps:

10

$$\frac{\partial \mathcal{L}}{\partial g^{[2]}} \Big|_K = \frac{1}{m} \sum_{e \in \{1, \dots, m\}}$$

$$\left( \frac{1 - \gamma^{(e)}}{1 - \gamma^{(e)}} - \frac{\gamma^{(e)}}{\gamma^{(e)}} \right)$$

$$\cdot \sigma(z^{[3](e)}) \cdot (1 - \sigma(z^{[3](e)}))$$

$$\left. \begin{array}{c} dz \\ \lambda^3 \end{array} \right\}$$

$$\lambda^3$$

This seems very transparent!  
if you need the transpose!

$$W^{[3]T} \cdot act^{[2]'}(z_1^{[2](e)})$$

$$p^{[2]} = w_{i,j}^{[2]}$$

$$q^{[2]} = b_j^{[2]}$$

$$\left\{ \begin{array}{c} p^{[2]} = w_{i,j}^{[2]} \\ q^{[2]} = b_j^{[2]} \end{array} \right.$$

$$\lambda^2$$

$$\lambda^2 W^{[3]T}$$

$$\left( \lambda^2 W^{[3]} \lambda^3 \right) dz^3$$

$$\int^{activation} *$$

$$act^{[2]'} \lambda^2$$

scalar

How does it look for layer 3?

(17)

$$\frac{\partial}{\partial g^{[1]} z^{[3]}(e)} = W^{[3]T} \cdot \frac{\partial}{\partial g^{[1]} z^{[2]}(e)}$$

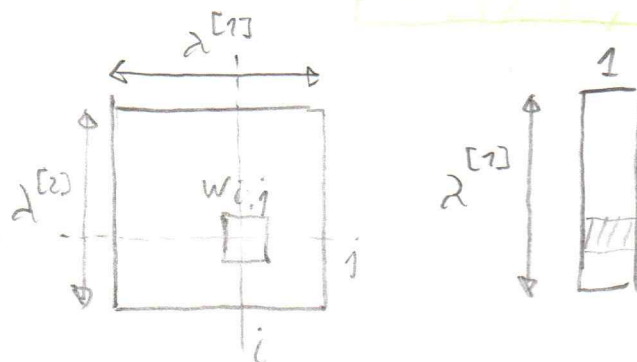
NB now it layer 1

$$\frac{\partial}{\partial g^{[1]} z^{[2]}(e)} = \text{act}^{[2]'} \left( z^{[2]}(e) \right) \cdot \frac{\partial}{\partial g^{[1]} z^{[1]}(e)}$$

Now we just need to work on the "red block" above:

$$\frac{\partial}{\partial g^{[1]} z^{[2]}(e)} = \frac{\partial}{\partial g^{[1]}} \left( W^{[2]T} \cdot a^{[1]}(e) + \underbrace{b^{[2]}}_{\text{resolves to } 0} \right)$$

$$= W^{[2]T} \cdot \frac{\partial}{\partial g^{[1]} a^{[1]}(e)}$$



Now we need to work on the "green block" above:

$$\frac{\partial}{\partial g^{[1]} a^{[1]}(e)}$$

For layer [1], we use an unspecified activation function  $\text{act}^{[1]}(-)$

$$= \underbrace{\text{act}^{[1]'} \left( z^{[1]}(e) \right)}_{\text{can be computed and cached during fwd computation}}$$

$$\cdot \frac{\partial}{\partial \phi^{[1]}} z^{[1]}(e)$$

can be computed and  
cached during fwd  
computation

Now we just need to work on the "red block" above:

$$\frac{\partial}{\partial \phi^{[1]}} z^{[1]}(e) = \frac{\partial}{\partial \phi^{[1]}} \left( W^{[1]T} \cdot a^{[0]}(e) + b^{[1]} \right)$$

Diagram illustrating the dimensions of the tensors in the equation above:

- $z^{[1]}(e)$  is a column vector of size  $1 \times \lambda^{[1]}$ .
- $W^{[1]T}$  is a matrix of size  $\lambda^{[0]} \times \lambda^{[1]}$ .
- $a^{[0]}(e)$  is a column vector of size  $1 \times \lambda^{[0]}$ .
- $b^{[1]}$  is a column vector of size  $1 \times \lambda^{[1]}$ .

► For  $\phi^{[1]} = w_{i,j}^{[1]}$ :

$$\frac{\partial}{\partial w_{i,j}^{[1]}} \left( W^{[1]T} \cdot a^{[0]}(e) + b^{[1]} \right) = a_i^{[0]}(e)$$

selects "a<sub>i</sub><sup>[0]</sup>(e)" yields  $\phi$

inside a column vector of  $\phi$ s

► For  $\phi^{[1]} = b_j^{[1]}$ :

$$\frac{\partial}{\partial b_j^{[1]}} \left( W^{[1]T} \cdot a^{[0]}(e) + b^{[1]} \right) = 1$$

yields  $\phi$  yields "1 at place j"

inside a column vector of  $\phi$ s

And thus

$$\left. \frac{\partial J}{\partial g^{[1]}} \right|_K = \frac{1}{m} \sum_{c \in \{1, \dots, m\}} \left( \frac{1 - y^{(c)}}{1 - \hat{y}^{(c)}} - \frac{y^{(c)}}{\hat{y}^{(c)}} \right)$$

$$\bullet \sigma(z^{[3](c)}) \cdot (1 - \sigma(z^{[3](c)}))$$

$$\bullet W^{[3]T} \cdot \text{act}^{[2]'}(z^{[2](c)})$$

$$\bullet W^{[2]T} \cdot \text{act}^{[1]'}(z^{[1](c)})$$

$$\bullet \left\{ g^{[1]} = w_{i,j}^{[1]} \right. ;$$

$$\left. g^{[1]} = b_j^{[1]} \right\} ; 1$$

□

