

# Backpropagation by hand

$$\text{Layer 2: } \frac{\partial \mathcal{L}}{\partial w_j^{(2)}} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial w_j^{(2)}} = (\hat{y} - y) \cdot \frac{\partial \hat{y}}{\partial w_j^{(2)}}$$

expanding the first summation:

$$\hat{y} = w_1^{(2)} \underbrace{\sum_i x_i w_{i1}^{(1)}}_{\frac{\partial \hat{y}}{\partial w_1^{(2)}}} + w_2^{(2)} \underbrace{\sum_i x_i w_{i2}^{(1)}}_{\frac{\partial \hat{y}}{\partial w_2^{(2)}}}$$

$$\text{So, } \frac{\partial \mathcal{L}}{\partial w_j^{(2)}} = (\hat{y} - y) \sum_i x_i w_{ij}^{(1)}$$

$$\text{Matrix form: } \hat{y} = X W^{(1)} W^{(2)}, \quad \frac{\partial \hat{y}}{\partial W^{(2)}} = X W^{(1)}$$

$$\therefore \frac{\partial \mathcal{L}}{\partial W^{(2)}} = (\hat{y} - y) X W^{(1)}$$

$$\text{Layer 1: } \frac{\partial \mathcal{L}}{\partial W^{(1)}} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial h} \frac{\partial h}{\partial W^{(1)}}, \text{ where } h = X W^{(1)}$$

$$= (\hat{y} - y) W^{(2)T} X$$

↑  
input to hidden layer

↑  
I missed this part!

$$\frac{\partial (X W^{(1)} W^{(2)})}{\partial (X W^{(1)})} = W^{(2)T}$$

with bias terms:

$$\hat{y} = (xw^{(1)} + b^{(1)})w^{(2)} + b^{(2)}$$

$$\hat{y} = C + 1b^{(2)}$$

$$\frac{\partial \hat{y}}{\partial b^{(2)}} = 1$$

$$\frac{\partial \mathcal{L}}{\partial b^{(2)}} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial b^{(2)}} = \frac{\partial \mathcal{L}}{\partial \hat{y}} = \hat{y} - y$$

with activations:

$$\hat{y} = f_2(\underbrace{f_1(xw^{(1)} + b^{(1)})}_{z_1}w^{(2)} + b^{(2)})$$

$z_2$

$$\frac{\partial \mathcal{L}}{\partial w^{(2)}} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial f_2} \frac{\partial f_2}{\partial z_2} \frac{\partial z_2}{\partial w^{(2)}}$$

①      ②      ③      ④

$$= (\hat{y} - y) \underbrace{(1)}_{\text{since } \hat{y} = f_2} \frac{\partial f_2}{\partial z_2} \underbrace{f_1(xw^{(1)} + b^{(1)})}_{\text{④}}$$

①      ②      ③      ④

assume both sigmoid

depends on function, e.g.  $f_2(z) = \sigma(z) = f_1(z)$

$$\frac{\partial f_2}{\partial z} = \sigma(z)(1 - \sigma(z))$$

store this

Layer 1:

$$\frac{\partial \mathcal{L}}{\partial w^{(1)}} = \left( \frac{\partial \mathcal{L}}{\partial f_2} \frac{\partial f_2}{\partial z_2} \right) \frac{\partial z_2}{\partial f_1} \frac{\partial f_1}{\partial z_1} \frac{\partial z_1}{\partial w^{(1)}}$$

$$= (\hat{y} - y) \frac{\partial f_2}{\partial z_2} w^{(2)\top} \frac{\partial f_1}{\partial z_1} x$$

$$= (\hat{y} - y) \sigma(z_2)(1 - \sigma(z_2)) w^{(2)\top} \sigma(z_1)(1 - \sigma(z_1)) x$$