

$$\text{Loss} = L = \frac{1}{2}(\hat{y} - y)^2 = \frac{1}{2}(0.751365 - 0.07)^2 \approx 0.27987$$

2) Backward pass - compute gradients

Compute the error signals (deltas) and then gradients for weights and biases

- Output Layer delta

$$\rightarrow \text{MSE} = L = \frac{1}{2}(\hat{y} - y)^2 \rightarrow \frac{\partial L}{\partial y} = \frac{\partial}{\partial \hat{y}} \left( \frac{1}{2}(\hat{y} - y)^2 \right)$$

$$\frac{\partial L}{\partial \hat{y}} = (\hat{y} - y) \rightarrow \text{derivative}$$

$$\text{derivative of sigmoid} = \sigma'(z) = \sigma(z)(1 - \sigma(z))$$

$$\sigma(z) = \frac{1}{1 + e^{-x}} \rightarrow \text{reciprocal} = \sigma(x) = (1 + e^{-x})^{-1}$$

apply the chain rule to the outer function

$$g(z) = (w^{(l+1)})^T g(z+1) \odot \sigma'(z^{(l)}) =$$

$$\rightarrow u \text{ with } u = 1 + e^{-x} \quad \frac{d}{dx} [(1 + e^{-x})^{-1}] = -(1 + e^{-x})^{-2} \cdot \frac{d}{dx} (1 + e^{-x})$$

$$\frac{d}{dx} (1 + e^{-x}) \rightarrow -e^{-x}$$

$$G'(x) = (1 + e^{-x})^{-2} \cdot (-e^{-x}) = \frac{e^{-x}}{(1 + e^{-x})^2} = \sigma'(x)$$

$$\frac{e^{-x}}{(1+e^{-x})^2} = \frac{1}{1+e^{-x}}$$

→ Solve for  $1-\sigma(x)$

$$1-\sigma(x) = 1 - \frac{1}{1+e^{-x}}$$

$$1-\sigma(x) = \frac{(1+e^{-x})}{(1+e^{-x})} - \frac{1}{(1+e^{-x})} = \frac{(1+e^{-x})-1}{1+e^{-x}}$$

$$\frac{e^{-x}}{1+e^{-x}} = 1-\sigma(x)$$

now → multiply  $\sigma(x)$  and  $(1-\sigma(x))$

$$\sigma(x)(1-\sigma(x)) = \frac{1}{1+e^{-x}} \cdot \frac{e^{-x}}{1+e^{-x}} = \frac{e^{-x}}{(1+e^{-x})^2}$$

$$\delta'(x) = \sigma(x)(1-\sigma(x)) = a(1-a)$$

if a function depend on another functions,  
we use the chain rule.  $\frac{dL}{da} = \frac{dL}{dz} \cdot \frac{da}{dz}$

$$\text{Delta} \rightarrow g = (\bar{y}-y) \cdot a(1-a)$$

$$g^{(2)}(1-a^{(2)}) = 0.751365(1-0.751365) \approx 0.186816$$

$$g^{(2)} = 0.7413565 \cdot 0.186816 \approx 0.138499$$

1

## Two-Layer Network

Input  $x = \begin{bmatrix} 0.05 \\ 0.10 \end{bmatrix}$ ,  $y = 0.01$ , Learning Rate = 0.5

Hidden layer size: 2 neurons

Output layer: 1 neuron

First-Layer weight & Bias

$$W^{(1)} = \begin{bmatrix} 0.15 & 0.20 \\ 0.25 & 0.30 \end{bmatrix}, b^{(1)} = \begin{bmatrix} 0.35 \\ 0.35 \end{bmatrix}$$

Second-Layer weights & Bias

$$W^{(2)} = \begin{bmatrix} 0.40 & 0.45 \end{bmatrix}, b^{(2)} = [0.60]$$

( $\rightarrow$  ~~prob~~)  $\rightarrow$  ~~hidden~~  $\rightarrow$  ~~output~~

Hidden Activation

$$z^{(1)} = (W^{(1)}x + b^{(1)}) \text{ weighted sum}$$

$$\text{neuron 1} \rightarrow = (0.15)(0.05) + (0.20)(0.10) + 0.35 = 0.3775$$

$$\text{neuron 2} \rightarrow = (0.25)(0.05) + (0.30)(0.10) + 0.35 = 0.3925$$

$$z^{(1)} = \begin{bmatrix} 0.3775 \\ 0.3925 \end{bmatrix} - z^{(1)}$$

Hidden Activation  $\sigma(z^{(1)}) = \frac{1}{1+e^{-x}}$

Neuron 1

$$\sigma(z_1^{(1)}) = 0.593269$$

Neuron 2

$$\sigma(z_2^{(1)}) = 0.596889$$

$$\sigma(z^{(1)}) = \begin{bmatrix} 0.593269 \\ 0.596889 \end{bmatrix} = a_2^{(1)}$$

$$W^{(2)} = [0.40 \ 0.45], b^{(2)} = [0.60]$$

$$\text{Output pre-activation: } z^{(2)} = w^{(2)}a^{(1)} + b^{(2)}$$

$$(0.40)(0.593269) + (0.45)(0.5916884) + 0.60$$

$$z^{(2)} = [1.1059054] + \dots$$

$$\text{Output-Activation} = \sigma(z^{(2)}) = \frac{1}{1+e^{-x}}$$

$$\hat{y} = \sigma(z^{(2)}) = 0.751365 \leftarrow a^{(2)}$$

$$\text{Loss Computation} \rightarrow L = \frac{1}{2}(\hat{y} - y)^2$$

$$L = 0.274811$$

(2)

## Backpropagation Deltas and Gradients

$$\sigma'(z) = \sigma(z)(1-\sigma(z)) \leftarrow \text{Sigmoid derivative}$$

$$\text{Derivative of Loss} \rightarrow \partial L / \partial a^{(2)} = \hat{y} - y$$

$$0.751365 - 0.01 = 0.741365$$

Activation derivative wrt pre-activation

~~$$\frac{\partial a^{(2)}}{\partial z^{(2)}} = \sigma'(z^{(2)}) = a^{(2)}(1-a^{(2)})$$~~

$$a^{(2)}(1-a^{(2)}) \approx 0.186816$$

$$\text{Chain rule multiplication: } g^{(2)} = \partial L / \partial a^{(2)} \cdot \partial a^{(2)} / \partial z^{(2)}$$

$$g^{(2)} = \partial L / \partial a^{(2)} \cdot \sigma'(z^{(2)}) = (y - \hat{y}) \cdot a^{(2)}(1-a^{(2)})$$

$$g^{(2)} = 0.741365(0.186816) = 0.138499$$

Gradients for  $W$  and  $b$

$$\frac{\partial L}{\partial w_1^{(2)}} = g^{(2)} \cdot (a^{(1)})^T \quad (\text{shape: } 7 \times 2)$$

$$\frac{\partial L}{\partial w_1^{(2)}} = 0.138499 (0.593269) = 0.082167 =$$

$$= \frac{\partial L}{\partial w_2^{(2)}} = 0.138499 (0.596884) = 0.082668$$

Bias gradient

$$\frac{\partial L}{\partial b^{(2)}} = g^{(2)} \approx 0.138499$$

Hidden layers

$$g^{(1)} = (w^{(2)})^T g^{(2)} \odot \sigma'(z^{(1)}) \quad [w^{(2)}]^T g^{(2)}$$

$$\text{For neuron 1: } 0.45 \cdot 0.138499 = 0.055400$$

$$\text{For neuron 2: } 0.45 \cdot 0.138499 = 0.062325$$

$$\sigma'(z^{(1)}) = a^{(1)}(1-a^{(1)})$$

$$= 0.593269(1-0.593269) = 0.241300$$

$$g^{(1)} = 0.593269(1-0.593269) = 0.241300$$

$$g_1^{(1)} = (0.055400) \cdot (0.241300) = 0.013368$$

$$g_2^{(1)} = (0.062325) \cdot (0.241300) = 0.0149974$$

Gradients for  $w^{(1)}$  and  $b^{(1)}$

$$\frac{\partial L}{\partial w^{(1)}} = g_1^{(1)} (x) \quad (\text{shape: } 2 \times 2)$$

$$= 0.013368(0.05) = 0.0006689$$

$$0.013368(0.10) = 0.0013368$$

$$g^{(1)} \cdot f(x) \rightarrow \text{inputs}$$

$$\begin{aligned} \text{Neuron } 2 &= 0.014974(0.05) = 0.0006684 + 0.007487 \\ &= 0.014974(0.10) = 0.0013368 + 0.0014974 \end{aligned}$$

- Bias gradient  $\frac{\partial L}{\partial b^{(1)}} = g^{(1)}$

$$\frac{\partial L}{\partial b^{(1)}} = g^{(1)} = \begin{bmatrix} 0.013368 \\ 0.014974 \end{bmatrix}$$

### 3) Parameter update (one gradient step)

- Gradient descent update:  $\theta \leftarrow \theta - \eta \cdot \nabla \theta$

- with  $\eta = 0.5$

- update  $W^{(2)}$  and  $b^{(2)}$

$$- W_{\text{new},11}^{(2)} = 0.40 - 0.5 \cdot 0.082167 \approx 0.358917$$

$$- W_{\text{new},21}^{(2)} = 0.45 - 0.5 \cdot 0.082668 \approx 0.408666$$

$$- b_{\text{new}}^{(2)} = 0.60 - 0.5 \cdot 0.138499 \approx 0.530751$$

$$- \text{So } W_{\text{new}}^{(2)} \approx [0.358917, 0.408666] \quad b$$

$$- b_{\text{new}}^{(2)} \approx 0.530751$$

~ Update  $W^{(1)}$  and  $b^{(1)}$

$$W_{\text{new},11}^{(1)} = 0.15 - 0.5 \cdot 0.0006689 \approx 0.14966$$

$$W_{\text{new},12}^{(1)} = 0.20 - 0.5 \cdot 0.0013368 \approx 0.199332$$

$$W_{\text{new},21}^{(1)} = 0.25 - 0.5 \cdot 0.007487 \approx 0.249626$$

$$W_{\text{new},22}^{(1)} = 0.30 - 0.5 \cdot 0.014974 \approx 0.299251$$

$$b_{\text{new},1}^{(1)} = 0.35 - 0.5 \cdot 0.013368 \approx 0.343316$$

$$b_{\text{new},2}^{(1)} = 0.35 - 0.5 \cdot 0.014974 \approx 0.342513$$

$$\therefore W_{\text{new}} = \begin{bmatrix} 0.14966 & 0.199332 \\ 0.249626 & 0.299251 \end{bmatrix} \quad b = \begin{bmatrix} 0.343316 \\ 0.342513 \end{bmatrix}$$

1 Forward pass to confirm  
the loss decreased

$$\text{New } z_1^{(1)} = 0.19966 \cdot 0.05 + 0.199332 \cdot 0.10 + 0.345316 \approx 0.37073$$

$$\text{New } z_2^{(1)} \approx 1.249626 \cdot 0.05 + 0.299251 \cdot 0.10 + 0.342513 \approx 0.38919$$

$$\text{New } a_1^{(1)} = \sigma(0.37073) \approx 0.591862$$

$$\text{New } a_2^{(1)} = \sigma(0.38919) \approx 0.595952$$

$$\text{New } z^{(2)} = 0.358917 + 0.591862 + 0.40866 \cdot 0.595952 + 0.530735 \\ \approx 0.98615$$

$$\text{New } y \approx \sigma(0.98615) \approx 0.727083$$

new loss  $\approx$

$$L_{\text{new}} = 0.5(0.72703 - 0.01)^2 \approx 0.2570$$

So the ~~gradient~~ <sup>Loss</sup> went down  
from  $\approx 0.274811$  to  $0.2570$