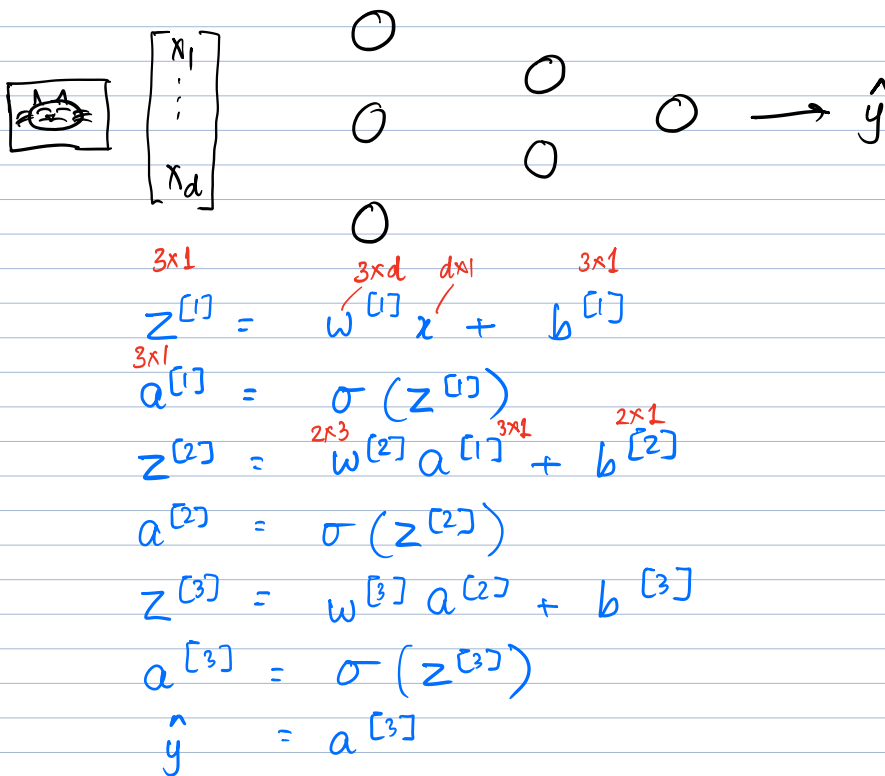


# Deep Learning

- ① Logistic Regression with a NN mindset
- ② Neural Networks
- ③ Backpropagation
- ④ Improving your NN



Optimizing  $w^{[1]}, w^{[2]}, w^{[3]}, b^{[1]}, b^{[2]}, b^{[3]}$

Loss (Cost) fn  $J(\hat{y}, y) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}^{(i)}$

$\mathcal{L}^{(i)} = -[y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log (1 - \hat{y}^{(i)})]$

$\frac{\partial J}{\partial w^{[2]}}$

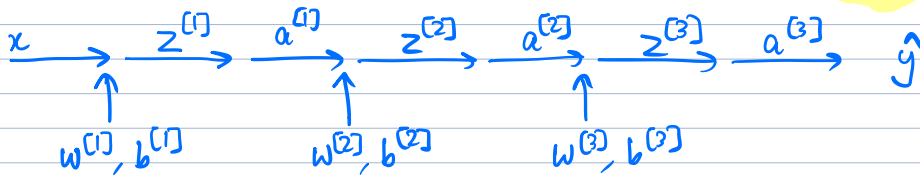
$\frac{\partial \mathcal{L}^{(i)}}{\partial w^{[2]}}$

Backward Propagation

$\forall l = 1 \dots 3$

$w^{[l]} = w^{[l]} - \alpha \frac{\partial J}{\partial w^{[l]}}$

$b^{[l]} = b^{[l]} - \alpha \frac{\partial J}{\partial b^{[l]}}$



$$\frac{\partial J}{\partial w^{[1]}} \text{ vs. } \frac{\partial J}{\partial w^{[3]}}$$

$$\frac{\partial J}{\partial w^{[3]}} = \frac{\partial J}{\partial a^{[3]}} \cdot \frac{\partial a^{[3]}}{\partial z^{[3]}} \cdot \frac{\partial z^{[3]}}{\partial w^{[3]}}$$

$$\frac{\partial J}{\partial w^{[2]}} = \frac{\partial J}{\partial z^{[3]}} \cdot \frac{\partial z^{[3]}}{\partial a^{[2]}} \cdot \frac{\partial a^{[2]}}{\partial z^{[2]}} \cdot \frac{\partial z^{[2]}}{\partial w^{[2]}}$$

$$\frac{\partial J}{\partial w^{[1]}} = \frac{\partial J}{\partial z^{[2]}} \cdot \frac{\partial z^{[2]}}{\partial a^{[1]}} \cdot \frac{\partial a^{[1]}}{\partial z^{[1]}} \cdot \frac{\partial z^{[1]}}{\partial w^{[1]}}$$

$$\cancel{\frac{\partial J}{\partial w^{[3]}}} \quad \frac{\partial \mathcal{L}^{(i)}}{\partial w^{[3]}} \quad \mathcal{L}^{(i)} = -[y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log (1 - \hat{y}^{(i)})]$$

$$\begin{aligned} \frac{\partial \mathcal{L}^{(i)}}{\partial w^{[3]}} &= -[y^{(i)} \frac{\partial}{\partial w^{[3]}} (\log \underbrace{\sigma(w^{[3]} a^{[2]} + b^{[3]})}_{a^{[3]}}) \\ &\quad + (1 - y^{(i)}) \frac{\partial}{\partial w^{[3]}} (\log (1 - \underbrace{\sigma(w^{[3]} a^{[2]} + b^{[3]})}_{a^{[3]}}))] \\ &= -[y^{(i)} \cdot \cancel{\frac{1}{a^{[3]}}} \cdot a^{[3]} (1 - a^{[3]}) \cdot a^{[2]T} \\ &\quad + (1 - y^{(i)}) \frac{1}{\cancel{1 - a^{[3]}}} \cdot (-1) \cdot a^{[3]} \cdot \cancel{(1 - a^{[3]})} \cdot a^{[2]T}] \end{aligned}$$

$$\frac{d}{dx} \log x = \frac{1}{x}$$

$$\frac{d}{dx} \sigma(x) = \sigma(x) (1 - \sigma(x))$$

$$\frac{\partial \log \sigma(\cdot)}{\partial w^{[3]}} = \frac{1}{\sigma(\cdot)} \cdot \frac{\partial \sigma(\cdot)}{\partial w^{[3]}}$$

$$\frac{\partial \sigma(w^{[3]} a^{[2]} + b^{[3]})}{\partial w^{[3]}} = a^{[3]} (1 - a^{[3]}) \cdot \frac{\partial (w^{[3]} a^{[2]} + b^{[3]})}{\partial w^{[3]}}$$

$$\stackrel{1 \times 2}{\frac{\partial (w^{[3]} a^{[2]} + b^{[3]})}{\partial w^{[3]}}} = \stackrel{1 \times 1}{a^{[2]T}}$$

$$\frac{\partial \mathcal{L}^{(i)}}{\partial w^{[3]}} = -[y^{(i)} \underbrace{(1 - a^{[3]})}_{-y^{(i)} a^{[3]}} \cdot a^{[2]T} - \underbrace{(1 - y^{(i)})}_{+ y^{(i)}} \cdot a^{[3]} \cdot a^{[2]T}]$$

$$= -(y^{(1)} - a^{(3)}) a^{(2)T}$$

$$\frac{\partial J}{\partial w^{(3)}} = -\frac{1}{n} \sum_{i=1}^n (y^{(1)} - a^{(3)}) a^{(2)T}$$

$$\frac{\partial \mathcal{L}}{\partial w^{(2)}} = \frac{\partial \mathcal{L}}{\partial a^{(3)}} \cdot \frac{\partial a^{(3)}}{\partial z^{(3)}} \cdot \frac{\partial z^{(3)}}{\partial a^{(2)}} \cdot \frac{\partial a^{(2)}}{\partial z^{(2)}} \cdot \frac{\partial z^{(2)}}{\partial w^{(2)}}$$

$$\frac{\partial \mathcal{L}}{\partial w^{(3)}} = \frac{\partial \mathcal{L}}{\partial z^{(3)}} \cdot \frac{\partial z^{(3)}}{\partial w^{(3)}} \quad a^{(2)T}$$

$$\frac{\partial \mathcal{L}^{(i)}}{\partial w^{(2)}} = (a^{(3)} - y^{(1)}) \cdot w^{(3)T} \cdot a^{(2)} (1 - a^{(2)}) \cdot a^{(1)T}$$

$$\begin{matrix} (2,3) & (1,1) & (2,1) & (2,1) & (1,3) \end{matrix}$$

$$= w^{(3)T} * a^{(2)} (1 - a^{(2)}) (a^{(3)} - y^{(1)}) \cdot a^{(1)T}$$

↑  
element  
wise  
product

$$\frac{\partial J}{\partial w^{(2)}} = \frac{1}{n} \sum_{i=1}^n \frac{\partial \mathcal{L}^{(i)}}{\partial w^{(2)}}$$

Try: compute  $\frac{\partial J}{\partial w^{(1)}}$   $\frac{\partial J}{\partial b^{(1)}}$

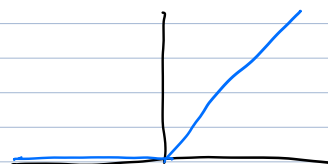
## Improving your NN

### ① Activation functions



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

$$\sigma'(z) = \sigma(z) (1 - \sigma(z))$$



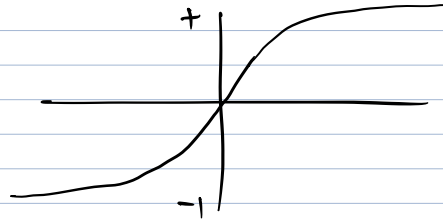
$$\text{ReLU}(z) = \begin{cases} 0 & \text{if } z \leq 0 \\ z & \text{if } z > 0 \end{cases}$$

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

replaces  $a^{(3)}(1 - a^{(3)})$   
 $a^{(2)}(1 - a^{(2)})$

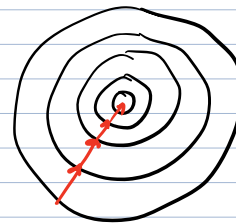
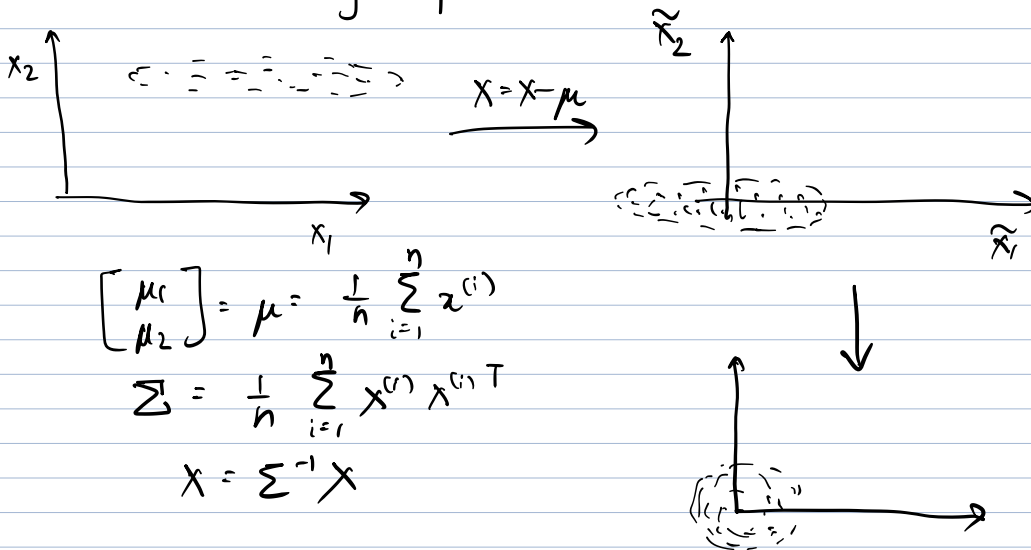
$$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$\tanh'(z) = 1 - \tanh(z)^2$$

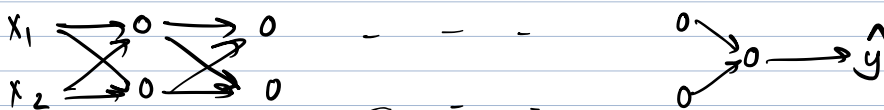


## ⑧ Initialization Methods

Normalizing input



Vanishing / Exploding gradients



Assume  $b=0$

activation  $f^n$  id:  $z \rightarrow z$

$$\hat{y} = w^{[L]} a^{[L-1]} = w^{[L]} w^{[L-1]} a^{[L-2]}$$

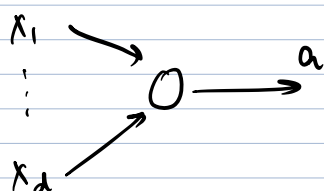
$$= w^{[L]} w^{[L-1]} \dots w^{[1]} x$$

$$W^{(t)} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix} \sim \begin{bmatrix} 1.5^L & 0 \\ 0 & 1.5^L \end{bmatrix}$$

$$\begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix} \sim \begin{bmatrix} 0.5^L & 0 \\ 0 & 0.5^L \end{bmatrix}$$

— avoid by initializing wts close to 1

Example w. one neuron



$$a = \sigma(z)$$

$$z = w_1 x_1 + \dots + w_d x_d$$

large  $d \rightarrow$  small  $w_i$

$$w_i \sim \frac{1}{d}$$

$$W_L = \text{np.random.randn}(\text{shape}) * \text{np.sqrt}\left(\frac{1}{n^{[L-1]}}\right)$$

↑  
for sigmoid

for ReLU : 2 instead of 1

Xavier Initialization

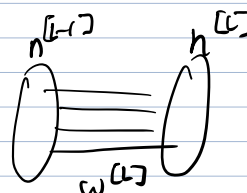
$$W^{[L]} \sim \sqrt{\frac{1}{n^{[L-1]}}} \quad \text{for tanh}$$

He Initialization

$$W^{[L]} \sim \sqrt{\frac{2}{n^{[L]} + n^{[L-1]}}}$$

↑  
backward  
prop.

↑  
forward  
prop



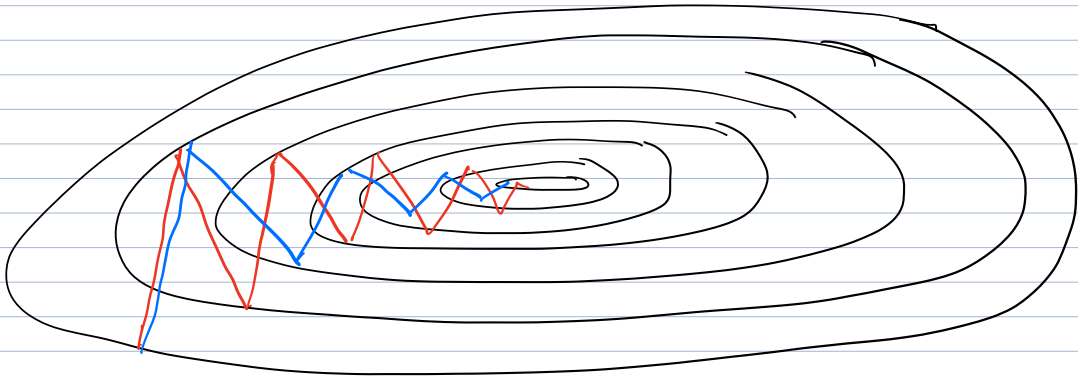
Optimization:

Gradient Descent

Mini batch gradient Descent



GD + Momentum



Momentum

$$W = W - \alpha \frac{\partial \mathcal{L}}{\partial W}$$

$$W = W - \alpha v$$

$$v = \beta \cdot v + (1 - \beta) \frac{\partial \mathcal{L}}{\partial W}$$

$v$ : momentum — weighted average of past updates