



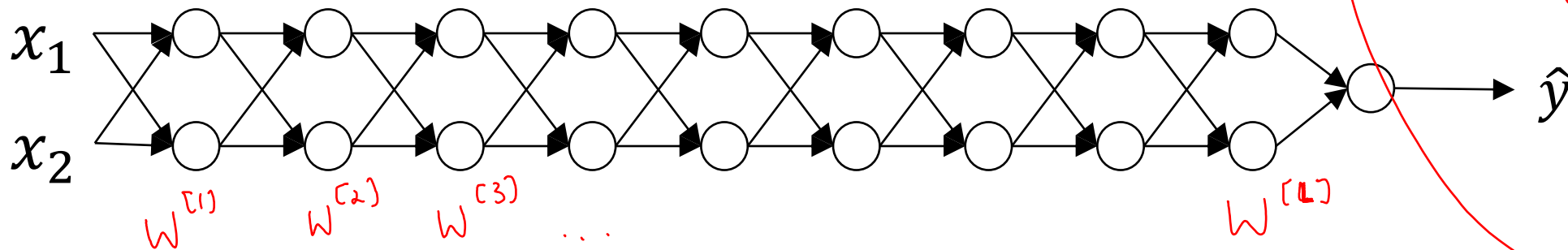
deeplearning.ai

Setting up your  
optimization problem

---

Vanishing/exploding  
gradients

# Vanishing/exploding gradients



This argument was for Activations!  
 $w^{(1)}$  is small  $\Rightarrow z$  is small  $\Rightarrow a$  is small, the same argument can be for their gradient  $dw$  is small  $\Rightarrow dz$  is small  $\Rightarrow da$  is small

$g(z) = z$  (Linear Activation Func)

$b^{(1)} = 0$  (All bias = 0)

$$\Rightarrow y = W^{(L)} \cdot W^{(L-1)} \cdot W^{(L-2)} \dots W^{(3)} \cdot W^{(2)} \cdot W^{(1)} \cdot X$$

$z^{(1)}$  (because  $b^{(1)} = 0$ )

Also  $g(z) = z = a$

$$\Rightarrow a^{(1)} = z^{(1)}$$

$$\Rightarrow W^{(1)} X = a^{(1)}$$

$$\Rightarrow W^{(2)} a^{(1)} = z^{(2)}, \text{ but } g(z^{(2)}) = a^{(2)} \Rightarrow W^{(2)} \cdot a^{(1)} = a^{(2)}$$

$$\Rightarrow W^{(L)} \cdot a^{(L-1)} = a^{(L)} = \hat{y}$$

dim for each  $w$

$$\hookrightarrow n^{(1)}, n^{(l-1)}$$

$$= 2 \times 2$$

Say we get values for each  $w$  as

$$W^{(2)} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}$$

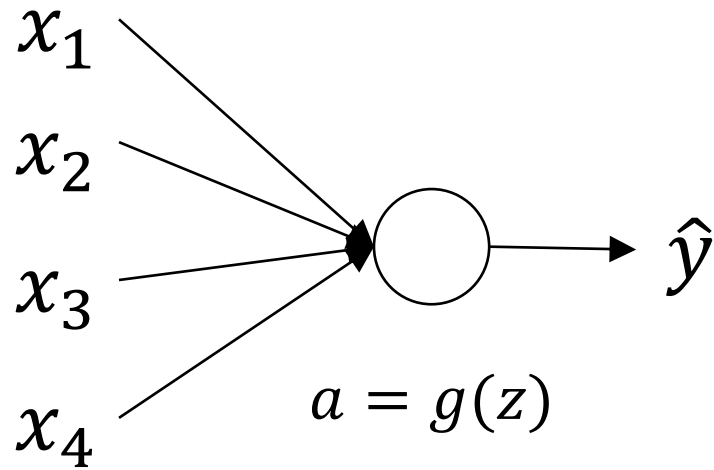
$$\Rightarrow \hat{y} = W^{(L)} \cdot \left[ \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix} \right]^{L-1} \cdot X \Rightarrow W^{(L)} \cdot \begin{bmatrix} 1.5^{L-1} & 0 \\ 0 & 1.5^{L-1} \end{bmatrix} \cdot X$$

to the power, i.e., matrix A, multiplied (L-1) times by itself

$\Rightarrow$  If "L" is a large number, say 30  $\Rightarrow 1.5^{(30)}$  is large  $\Rightarrow$  exploding gradient

Conversely If  $W^{(2)} = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}$  then we have  $(0.5)^{30} \Rightarrow$  vanishing gradient

# Single neuron example



(weight Initialization for deep neural networks)  
 ↳ partial solution to vanishing/exploding gradients is careful initialization

This doesn't solve vanishing/exploding gradient, but ↓ it significantly  
 ⇒ So as long as we normalize i/p's Initialize  $w$ 's like this, we can solve this problem

other variants

If using a tanh Activation use

$$\sqrt{\frac{1}{n^{(l-1)}}} \quad \text{or}$$

$$\sqrt{\frac{2}{n^{(l-1)} + n^{(l)}}} \quad \text{instead}$$

$$\text{of } \sqrt{\frac{1}{n^{(l-1)}}}$$

$z = w_1x_1 + w_2x_2 + \dots + w_nx_n + b$   
 If  $n$  is large, i.e., # features is large ⇒ you want all  $w_i$  to be small

- one way to do this

✗  $\text{Var}(w_i) = \frac{1}{n}$ , can also make  $\text{Var}(w_i) = \frac{2}{n}$  [same thing, prop to  $1/n$ ]

If  $\rightarrow W^{(L)} = \text{np.random.randn}(\text{shape}(w)) * \boxed{\text{np.sqrt}\left(\frac{2}{n^{(l-1)}}\right)}$  ⇒ then  $\text{Var}(w^{(l)}) = \frac{2}{n}$

↳ # features feeding into layer "l"

Note  
 If you are using ReLU Activation then setting  $\text{Var}(w_i) = \frac{2}{n}$  is better

May be good to make  $\text{Var}(w)$  a hyper param to change based on Activation func