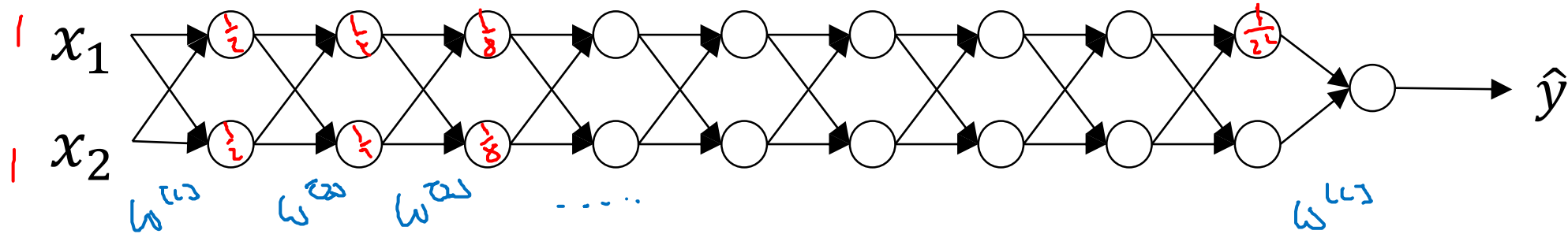deeplearning.ai

# Setting up your optimization problem

## Vanishing/exploding gradients

Your derivatives or your slopes can sometimes get either very very big, or very very small, maybe exponentially small, and this makes training difficult.
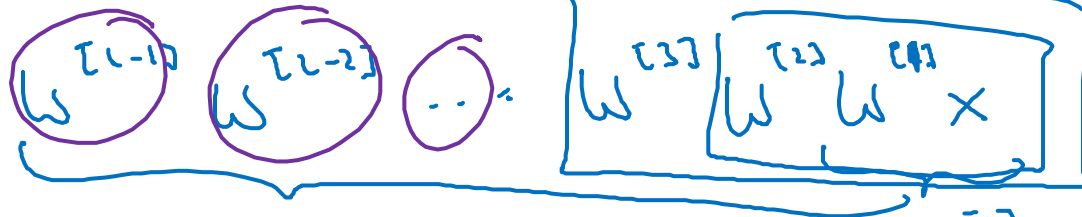
# Vanishing/exploding gradients

$L = 150$

This makes training difficult, especially if your gradients are exponentially smaller than L. Then gradient descent will take tiny little steps. It will take a long time for gradient descent to learn anything.



$1$ $x_1$

$1$ $x_2$

$\frac{1}{2}$ $\frac{1}{4}$ $\frac{1}{8}$ $\frac{1}{2^L}$

$\frac{1}{2}$ $\frac{1}{4}$ $\frac{1}{8}$

$W^{[1]}$ $W^{[2]}$ $W^{[3]}$ $\cdots$ $W^{[L]}$

$\hat{y}$

$g(z) = z.$   $b^{[l]} = 0.$

$\hat{y} = W^{[L]} \; W^{[L-1]} \; W^{[L-2]} \; \cdots \; W^{[3]} \; W^{[2]} \; W^{[1]} \; x$   $a^{[3]}$

$1.5^L$

$0.5^L$

Identity matrices

Increasing/Decreasing exponentially.

$W^{[l]} > I$

$W^{[l]} < I$   $\begin{bmatrix} 0.9 & \\ & 0.9 \end{bmatrix}$

$W^{[l]} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}$   $\begin{matrix} 0.5 \\ 0.5 \end{matrix}$

$z^{[1]} = W^{[1]} x$

$a^{[1]} = g(z^{[1]}) = z^{[1]}$

$a^{[2]} = g(z^{[2]}) = g(W^{[2]} a^{[1]})$

$\hat{y} = W^{[L]} \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}^{L-1} x$   $\begin{matrix} 0.5 \\ 0.9 \end{matrix}$

$1.5^{L-1} x$

$0.5^{L-1} x$

Andrew Ng

# Single neuron example

$$a^{[1]}$$

$x_1$

$x_2$

$x_3$

$x_4$

$$\hat{y}$$

$$W^{[1]}$$

$$a = g(z)$$

$$z = W_1 X_1 + W_2 X_2 + \cdots + W_n X_n \not{+b}$$

Large $n \rightarrow$ Smaller $W_i$

n #of features

$$Var(W_i) = \frac{1}{n} \frac{2}{n}$$

$$W^{[\ell]} = np.random.randn(shape..) * np.sqrt\left(\frac{2}{n^{[\ell-1]}}\right)$$

ReLU

$$g^{[\ell]}(z) = ReLU(z)$$

Other variants:

tanh

$$\frac{1}{n^{[\ell-1]}}$$

Xavier initialization

$$\frac{2}{n^{[\ell-1]} + n^{[\ell]}}$$

Andrew Ng