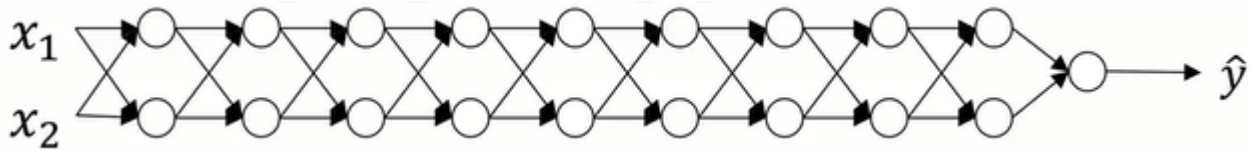


Vanishing/exploding gradients

Example



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

$$y = W^{[L]} W^{[L-1]} \dots W^{[2]} W^{[1]} x$$

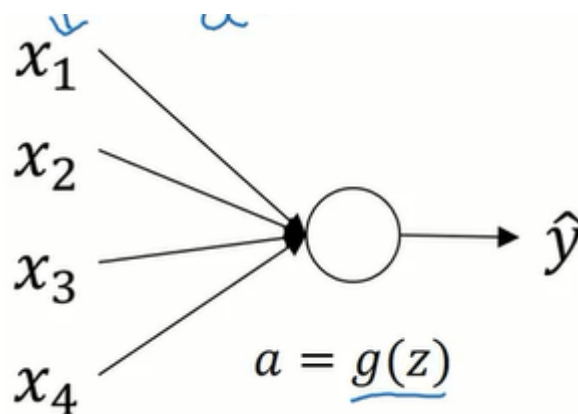
$$W^{[l]} = \begin{bmatrix} a & 0 \\ 0 & a \end{bmatrix}$$

$$y = W^{[L]} \begin{bmatrix} a & 0 \\ 0 & a \end{bmatrix}^{(L-1)} x = W^{[L]} \begin{bmatrix} a^{(L-1)} & 0 \\ 0 & a^{(L-1)} \end{bmatrix} x$$

$$\begin{cases} \text{exploding if } a > 1 \\ \text{vanishing if } a < 1 \end{cases}$$

Weight initialization for deep networks

•



$$z = w_1 x_1 + w_2 x_2 + \dots + w_n x_n$$

▪ want z not to explode, larger $n \rightarrow$ smaller w_i

$$\text{Var}(w_i) = \begin{cases} \frac{2}{n}, & \text{when } g(z) \text{ is ReLU} \\ \frac{1}{n}, & \text{when } g(z) \text{ is tanh} \end{cases}$$

$$W^{[l]} = \begin{cases} \text{np.random.randn(shape)} * \text{np.sqrt}(\frac{2}{n^{[l-1]}}), & \text{when } g^{[l]}(z) = \text{ReLU}(z) \\ \text{np.random.randn(shape)} * \text{np.sqrt}(\frac{1}{n^{[l-1]}}), & \text{when } g^{[l]}(z) = \text{tanh}(z) \end{cases}$$