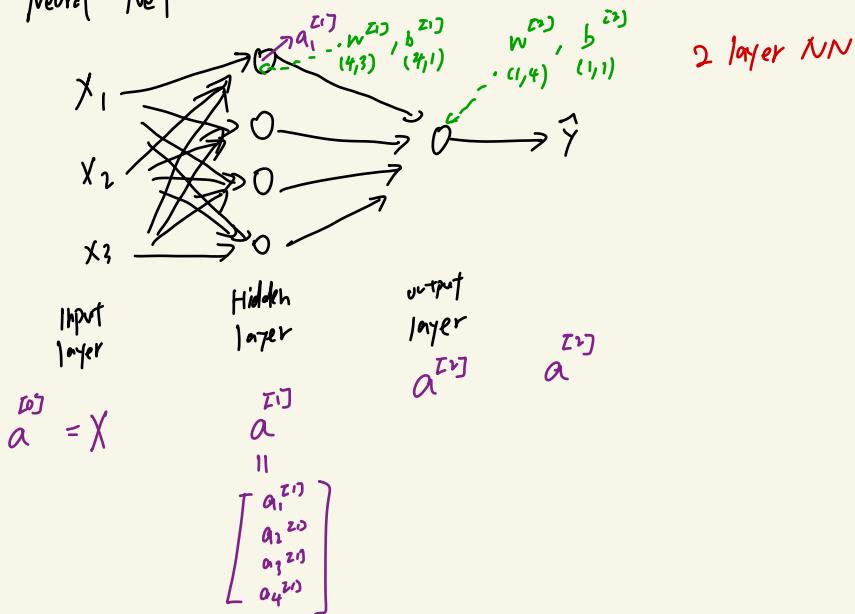


$[1] \leftarrow \text{layer}$
X

Neural Net



Hidden layer node

$$z_1^{[1]} = w_1^{[1]} x + b_1^{[1]} \rightarrow a_1^{[1]} = \sigma(z_1^{[1]})$$

$$\begin{aligned} X &\rightarrow a^{[2]} = \hat{y} \\ x^{(1)} &\rightarrow a^{[2](1)} = \hat{y}^{(1)} \\ &\vdots \\ &a^{[2](i)} \end{aligned}$$

for i in n

$$\begin{aligned} z^{[1](i)} &= w^{[1]} x^{(i)} + b^{[1]} \\ a^{[2](i)} &= \sigma(z^{[1](i)}) \end{aligned}$$

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & x^{(3)} & \dots & x^{(n)} \end{bmatrix}$$

$$A^{[1]} = \begin{bmatrix} a^{1} & a^{[1](2)} & \dots & a^{[1](m)} \end{bmatrix} \quad \downarrow \# \text{ hidden unit}$$

training sample \rightarrow

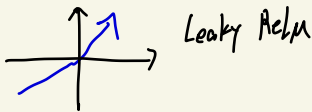
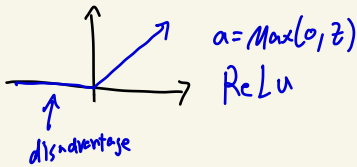
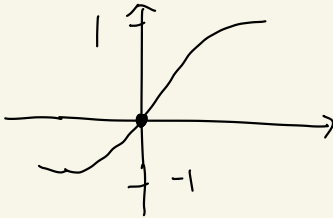
$$W' = \begin{bmatrix} \equiv \end{bmatrix} \quad W'X' = \begin{bmatrix} \vdots \end{bmatrix} \quad W'X'^2 = \begin{bmatrix} \vdots \end{bmatrix} \quad W'X'^3 = \begin{bmatrix} \vdots \end{bmatrix}$$

$$W' \begin{bmatrix} | & | & | & \dots & | \\ x^1 & x^2 & x^3 & \dots & x^n \\ | & | & | & \dots & | \end{bmatrix} = \begin{bmatrix} \bullet & \bullet & \bullet & \dots & \bullet \\ \vdots & \vdots & \vdots & \dots & \vdots \\ \bullet & \bullet & \bullet & \dots & \bullet \end{bmatrix} = \begin{bmatrix} z^{[1]} & z^{[2]} & z^{[3]} \\ | & | & | \\ +b^{[1]} & +b^{[2]} & +b^{[3]} \end{bmatrix} = \frac{z}{2}$$

X

Activation function

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



Why use activation function?

$a^{[2]}$ just linear output

If $0 \leq \hat{y} \leq 1$, sigmoid

tanh is always better

If $-1 \leq \hat{y} \leq 1$, $\tanh(z)$

$$a = g(z) = \frac{1}{1 + e^{-z}}$$

$$g'(z) = a(1-a)$$

$$g(z) = \tanh(z)$$

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

ReLU

$$g(z) = \max(0, z)$$

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

Leaky ReLU

$$g(z) = \max(0.01z, z)$$

$$g'(z) = \begin{cases} 0.01 & \text{if } z < 0 \\ 1 & \text{if } z \geq 1 \end{cases}$$

Gradient Descent for Neural Net

parameters: $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}$

\downarrow \downarrow \downarrow
 $(n^{[1]}, n^{[2]})$ $(n^{[1]}, n^{[2]})$ $(n^{[2]}, 1)$
 \downarrow
 $(n^{[1]}, 1)$

cost function $J = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(\hat{y}, y)$

\uparrow
 $a^{[2]}$

$h^x = h^{[2]}, h^{[1]}, n^{[2]} = 1$

Random init

$w^{[1]} = \text{np.random.randn}(n^{[1]}, n^{[2]}) \times 0.01$

$b^{[1]} = \text{np.zeros}(n^{[1]})$

Gradient Descent:

init random

Repeat \leq

predict $(\hat{y}, 1 \dots m)$

$dw = \frac{dJ}{dw}$ $db = \frac{dJ}{db}$

$w^{[1]} = w^{[1]} - \alpha dw^{[1]}$

$b^{[1]} = b^{[1]} - \alpha db^{[1]}$

$w^{[2]} = w^{[2]} - \alpha dw^{[2]}$

Forward

$z^{[1]} = w^{[1]} X + b^{[1]}$

$A^{[1]} = \sigma^{[1]}(z^{[1]})$

$z^{[2]} = w^{[2]} A^{[1]} + b^{[2]}$

$A^{[2]} = \sigma^{[2]}(z^{[2]})$

Back

$dz^{[2]} = A^{[2]} - Y$

$Y = [Y^{(1)} Y^{(2)} \dots Y^{(m)}]$

$dw^{[2]} = \frac{1}{n} dz^{[2]} A^{[1]T}$

$db^{[2]} = \frac{1}{n} \text{np.sum}(dz^{[2]}, \text{axis}=1, \text{keepdim}=True)$

$dz^{[1]} = w^{[2]T} dz^{[2]} + \sigma^{[1]'}(z^{[1]}) dz^{[2]}$ stop \downarrow ch, $\rightarrow (n, 1)$

\nwarrow \uparrow
 $(n^{[1]}, n)$ $(n^{[2]}, m)$
 element wise product

$dw^{[1]} = \frac{1}{n} dz^{[1]} X^T$

$db^{[1]} = \frac{1}{n} \text{np.sum}(dz^{[1]}, \text{axis}=1, \text{keepdim}=True)$