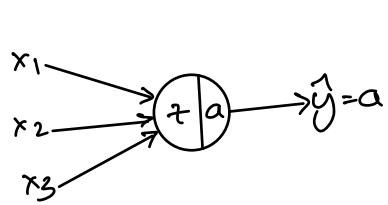
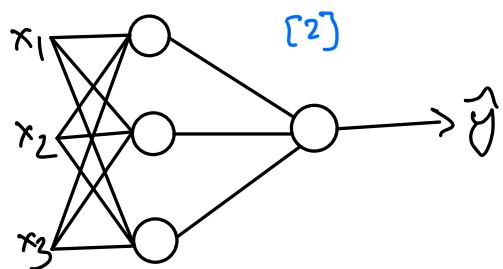


defn: One hidden layer Network:-



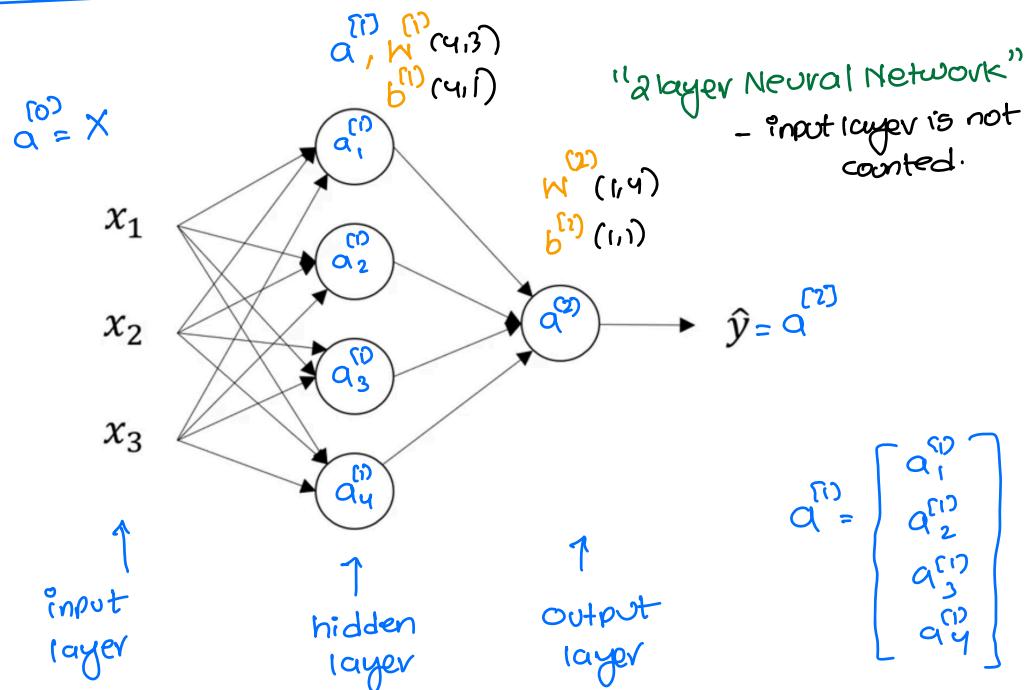
$$x \begin{matrix} \omega \\ b \end{matrix} \rightarrow z = w^T x + b \rightarrow a = \sigma(z) \rightarrow \hat{y} = \sigma(a)$$

[1]

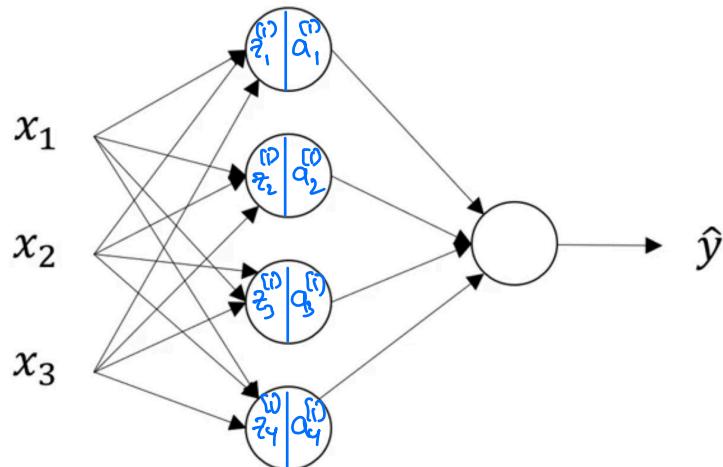
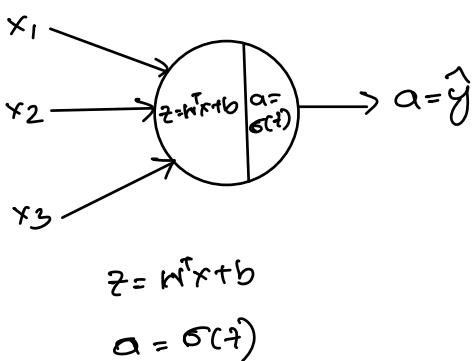


$$x \begin{matrix} d\omega^{[1]} \\ d\theta^{[1]} \\ dw^{[1]} \\ db^{[1]} \end{matrix} \rightarrow z^{[1]} = w^{[1]} x + \theta^{[1]} \rightarrow a^{[1]} = \sigma(z^{[1]}) \rightarrow \begin{matrix} d\omega^{[2]} \\ d\theta^{[2]} \\ dw^{[2]} \\ db^{[2]} \end{matrix} \rightarrow z^{[2]} = w^{[2]} a^{[1]} + \theta^{[2]} \rightarrow a^{[2]} = \sigma(z^{[2]}) \rightarrow \hat{y} = \sigma(a^{[2]})$$

Lec 2: Neural Network Representation



Lec 3: Computing a NN's output!



$$z_1^{[1]} = \omega_1^{[1] T} x + b_1^{[1]} \quad a_1^{[1]} = \sigma(z_1^{[1]})$$

$$z_2^{[1]} = \omega_2^{[1] T} x + b_2^{[1]} \quad a_2^{[1]} = \sigma(z_2^{[1]})$$

$$z_3^{[1]} = \omega_3^{[1] T} x + b_3^{[1]} \quad a_3^{[1]} = \sigma(z_3^{[1]})$$

$$z_4^{[1]} = \omega_4^{[1] T} x + b_4^{[1]} \quad a_4^{[1]} = \sigma(z_4^{[1]})$$

$$\begin{aligned}
 z_1^{(1)} &= w_1^{(1)\top} x + b_1^{(1)} & a_1^{(1)} &= \sigma(z_1^{(1)}) \\
 z_2^{(1)} &= w_2^{(1)\top} x + b_2^{(1)} & a_2^{(1)} &= \sigma(z_2^{(1)}) \\
 z_3^{(1)} &= w_3^{(1)\top} x + b_3^{(1)} & a_3^{(1)} &= \sigma(z_3^{(1)}) \\
 z_4^{(1)} &= w_4^{(1)\top} x + b_4^{(1)} & a_4^{(1)} &= \sigma(z_4^{(1)})
 \end{aligned}$$

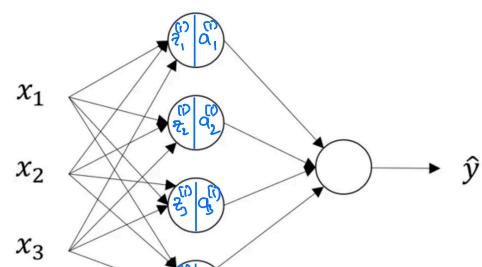
vectorize

$$\left[\begin{array}{c} -w_1^{(1)\top} \\ -w_2^{(1)\top} \\ -w_3^{(1)\top} \\ -w_4^{(1)\top} \end{array} \right] \left[\begin{array}{c} x_1 \\ x_2 \\ x_3 \end{array} \right] + \left[\begin{array}{c} b_1^{(1)} \\ b_2^{(1)} \\ b_3^{(1)} \\ b_4^{(1)} \end{array} \right] = \left[\begin{array}{c} z_1^{(1)} \\ z_2^{(1)} \\ z_3^{(1)} \\ z_4^{(1)} \end{array} \right] = z^{(1)}$$

$\begin{matrix} 4 \times 3 \\ N \end{matrix}$ x $b^{(1)}$ $z^{(1)}$

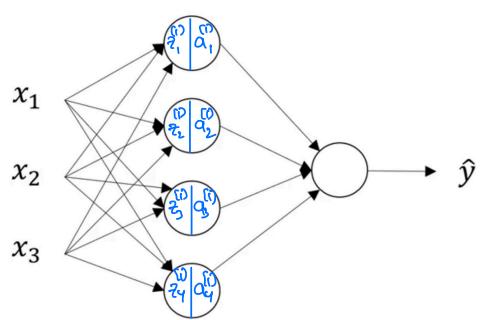
$$a^{(1)} = \begin{bmatrix} a_1^{(1)} \\ a_2^{(1)} \\ a_3^{(1)} \\ a_4^{(1)} \end{bmatrix} = \sigma(z^{(1)})$$

Final equations: $a^{(0)} = x_{(3,1)}$



$$\begin{aligned}
 z^{(0)} &= w^{(0)\top} a^{(0)} + b^{(0)} \\
 a^{(0)} &= \sigma(z^{(0)}) \\
 z^{(1)} &= w^{(1)\top} a^{(0)} + b^{(1)} \\
 a^{(1)} &= \sigma(z^{(1)})
 \end{aligned}$$

decr: Vectorizing across multiple examples:-



$$\begin{aligned} a^{(0)} &= x \\ z^{(1)} &= W^{(1)} a^{(0)} + b^{(1)} \\ a^{(1)} &= \sigma(z^{(1)}) \\ z^{(2)} &= W^{(2)} a^{(1)} + b^{(2)} \\ a^{(2)} &= \sigma(z^{(2)}) \end{aligned}$$

for single training example.

$$\begin{aligned} x &\longrightarrow a^{(1)} = \hat{y} \\ x^{(1)} &\longrightarrow a^{(1)(1)} = \hat{y}^{(1)} \\ x^{(1)} &\longrightarrow a^{(1)(2)} = \hat{y}^{(1)} \\ x &\longrightarrow a^{(2)} = \hat{y}^{(2)} \\ &\vdots \\ x^{(m)} &\longrightarrow a^{(m)} = \hat{y}^{(m)} \end{aligned}$$

$a^{(i,j)}$
i-th layer
j-th example.

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix}_{(n_x, m)}$$

for $i = 1$ to m :

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

Vectorized:-

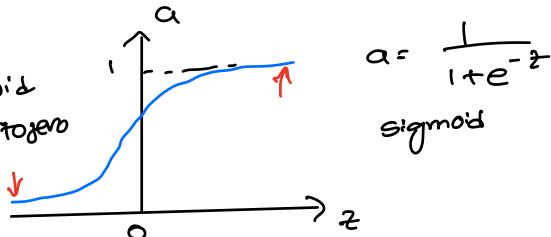
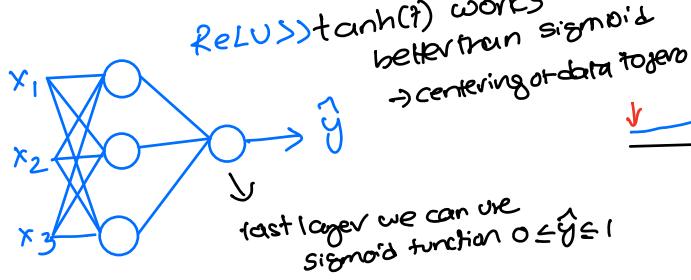
$$\begin{aligned} 1 \quad z^{(1)} &= W^{(1)} A^{(0)} + b^{(1)} \\ A^{(1)} &= \sigma(z^{(1)}) \end{aligned}$$

$$\begin{aligned} 2 \quad z^{(2)} &= W^{(2)} A^{(1)} + b^{(2)} \\ A^{(2)} &= \sigma(z^{(2)}) \end{aligned}$$

$$A^{(2)} = \begin{bmatrix} a^{(2)(1)} & a^{(2)(2)} & \dots & a^{(2)(m)} \end{bmatrix}$$

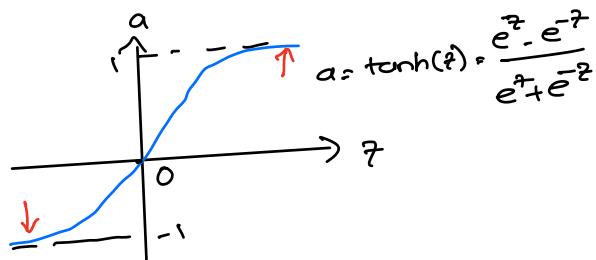
$$A^{(1)} = \begin{bmatrix} a^{(1)(1)} & a^{(1)(2)} & \dots & a^{(1)(m)} \\ \vdots & & & \\ a^{(1)(1)} & a^{(1)(2)} & \dots & a^{(1)(m)} \end{bmatrix}_{(q, m)}$$

defcs:- Activation Function:-

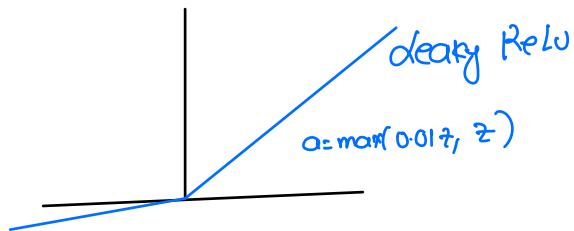
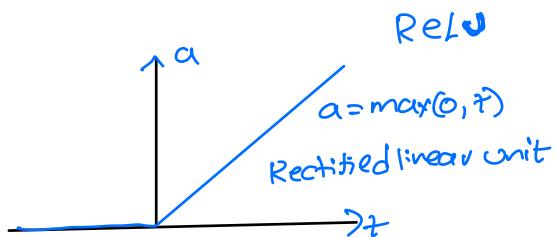


Given x

$$\begin{aligned} z^{[1]} &= W^{[1]}x + b^{[1]} \\ a^{[1]} &= \sigma(z^{[1]}) \quad g(z^{[1]}) \text{ tanh} \\ z^{[2]} &= W^{[2]}a^{[1]} + b^{[2]} \\ a^{[2]} &= \sigma(z^{[2]}) \quad g(z^{[2]}) \text{ sigmoid} \end{aligned}$$

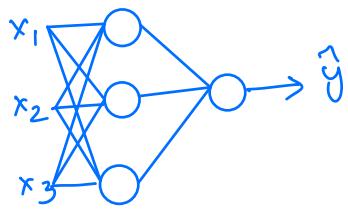


da is zero \rightarrow gradient descent becomes slow.



Ques 6:- Why do you need non linear activation function?

assume that there is no activation function.



Given x :

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

$$a^{[1]} = f^{[1]} \quad (\text{linear activation})$$

$$z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$$

$$a^{[2]} = z^{[2]}$$

$$a^{[1]} = W^{[1]} x + b^{[1]}$$

$$a^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$$

$$= W^{[2]} [W^{[1]} x + b^{[1]}] + b^{[2]}$$

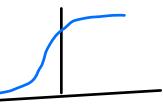
$$a^{[2]} = [W^{[2]} W^{[1]}] x + W^{[2]} b^{[1]} + b^{[2]}$$

$$a^{[2]} = W x + b \quad (\text{linear output})$$

if $y \in \mathbb{R}$, we can use $g(z) = z$ linear activation, such as
in output layer in regression.

lect 7: Derivatives of Activation Functions:-

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



$$\frac{dg(z)}{dz} = \frac{-1(e^{-z})}{(1+e^{-z})^2} = \frac{e^{-z}}{1+e^{-z}} \cdot \frac{1}{1+e^{-z}} = (1-g(z))g(z)$$

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

$$z=10 \quad g(z) \approx 1 \quad g'(z) \approx 0$$

$$z=-10 \quad g(z) \approx 0 \quad g'(z) \approx 0$$

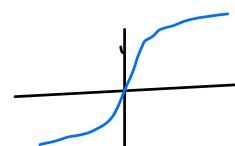
$$z=0 \quad g(z)=0.5 \quad g'(z)=0.25$$

Tanh activation :-

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

$$g'(z) = \frac{(e^z + e^{-z})(e^z - e^{-z}) - (e^z - e^{-z})(e^z + e^{-z})}{(e^z + e^{-z})^2}$$

$$= 1 - \left(\frac{e^z - e^{-z}}{e^z + e^{-z}} \right)^2$$



$$z=0 \quad \tanh(z) \approx 0 \quad g'(z) \approx 0$$

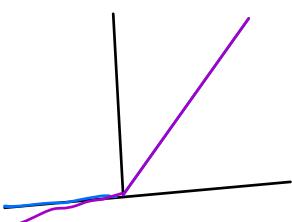
$$z=-10 \quad \tanh(z) \approx -1 \quad g'(z) \approx 0$$

$$z=10 \quad \tanh(z) \approx 1 \quad g'(z) \approx 0$$

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

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

ReLU, Leaky



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

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

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

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

dec8:- Gradient Descent for single hidden layer:-

parameters: $\omega^{(1)}, b^{(1)}, \omega^{(2)}, b^{(2)}$ $n_x = n^{(0)}$ - input features

$W^{(1)} = (n^{(0)}, n^{(1)})$ $b^{(1)} = (n^{(1)})$ $n^{(1)}$ - hidden units

$W^{(2)} = (n^{(1)}, n^{(2)})$ $b^{(2)} = (n^{(2)})$ $n^{(2)} = 1$ - output units

Cost function: $J(\omega^{(1)}, b^{(1)}, \omega^{(2)}, b^{(2)}) = \frac{1}{m} \sum_{i=1}^m \alpha(y_i, \hat{y}_i) \quad \hat{y} = a^{(2)}$

Initialize $\omega, b, \omega^{(1)}, b^{(1)}, \omega^{(2)}, b^{(2)}$

Gradient descent:

Repeat \leftarrow

compute prediction ($\hat{y}^{(i)}, i=1, \dots, m$)

$d\omega^{(1)}, db^{(1)}, d\omega^{(2)}, db^{(2)}$ $\frac{\partial J}{\partial \text{var}}$

update $\omega^{(1)} = \omega^{(1)} - \alpha d\omega^{(1)}$

$b^{(1)} = b^{(1)} - \alpha db^{(1)}$

$\omega^{(2)} = \omega^{(2)} - \alpha d\omega^{(2)}$

$b^{(2)} = b^{(2)} - \alpha db^{(2)}$

\downarrow

Equations: α eqn

Forward propagation:

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

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

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

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

back propagation:-

$$dZ^{[2]} = A^{[2]} - Y \quad Y = [y^{(1)}, \dots, y^{(m)}]_{(1, m)}$$

$$dW^{[2]} = \frac{1}{m} dZ^{[2]} A^{[1]T}$$

$$db^{[2]} = \frac{1}{m} np.sum(dZ^{[2]}, axis=1, keepdims=True) \quad (n^{[2]}, 1)$$

$$dZ^{[1]} = \underbrace{W^{[2]T} dZ^{[2]}}_{(n^{[1]}, m)} * \underbrace{g^{[1]}'(z^{[1]})}_{(n^{[1]}, m)} \quad (\text{elementwise})$$

$$dW^{[1]} = \frac{1}{m} dZ^{[1]} X^T$$

$$db^{[1]} = \frac{1}{m} np.sum(dZ^{[1]}, axis=1, keepdims=True) \quad (n^{[1]}, 1)$$

Summary of gradient descent

$$dz^{[2]} = a^{[2]} - y$$

$$dW^{[2]} = dz^{[2]} a^{[1]T}$$

$$db^{[2]} = dz^{[2]}$$

$$dz^{[1]} = W^{[2]T} dz^{[2]} * g^{[1]}'(z^{[1]}) \quad (n^{[1]}, 1)$$

$$dW^{[1]} = dz^{[1]} X^T$$

$$db^{[1]} = dz^{[1]}$$

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

$$dW^{[2]} = \frac{1}{m} dZ^{[2]} A^{[1]T}$$

$$db^{[2]} = \frac{1}{m} np.sum(dZ^{[2]}, axis=1, keepdims=True)$$

$$dZ^{[1]} = \underbrace{W^{[2]T} dZ^{[2]}}_{(n^{[1]}, m)} * \underbrace{g^{[1]}'(z^{[1]})}_{(n^{[1]}, m)} \quad (\text{elementwise product})$$

$$dW^{[1]} = \frac{1}{m} dZ^{[1]} X^T$$

$$db^{[1]} = \frac{1}{m} np.sum(dZ^{[1]}, axis=1, keepdims=True)$$

Summary of gradient descent (W3AI)

$$dz^{[2]} = a^{[2]} - y$$

$$dW^{[2]} = dz^{[2]} a^{[1]T}$$

$$db^{[2]} = dz^{[2]}$$

$$dz^{[1]} = W^{[2]T} dz^{[2]} * g^{[1]'}(z^{[1]})$$

$$dW^{[1]} = dz^{[1]} x^T$$

$$db^{[1]} = dz^{[1]}$$

$$\begin{matrix} \{0\} & \{1\} \\ \textcircled{0} & \textcircled{0} \\ \textcircled{0} & \textcircled{0} \\ \textcircled{0} & \textcircled{0} \\ \textcircled{2} & \textcircled{4} \end{matrix}$$

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

$$dW^{[2]} = \frac{1}{m} dZ^{[2]} A^{[1]T}$$

$$db^{[2]} = \frac{1}{m} np.sum(dZ^{[2]}, axis=1, keepdims=True)$$

$$dZ^{[1]} = \underbrace{W^{[2]T} dZ^{[2]}}_{(n^{[2]}, m)} * \underbrace{g^{[1]'}(Z^{[1]})}_{\substack{\text{elementwise product} \\ (n^{[1]}, m)}}$$

$$dW^{[1]} = \frac{1}{m} dZ^{[1]} X^T$$

$$db^{[1]} = \frac{1}{m} np.sum(dZ^{[1]}, axis=1, keepdims=True)$$

$$\begin{matrix} \{2\} \\ \textcircled{0} \rightarrow \hat{y} \\ \textcircled{1} \end{matrix} \quad \begin{matrix} W_1 = (4, 2) \\ b_1 = (4, 1) \end{matrix} \quad \begin{matrix} W_2 = (1, 4) \\ b_2 = (1, 1) \end{matrix} \quad \begin{matrix} X = (2, m) \\ Y = (1, m) \end{matrix}$$

Forward propagation:

$$z^{(1)} = W^{(1)} x + b^{(1)}$$

$$A^{(1)} = g^{(1)}(z^{(1)})$$

$$z^{(2)} = W^{(2)} A^{(1)} + b^{(2)}$$

$$A^{(2)} = g^{(2)}(z^{(2)})$$

$$z^{(1)} = (4, m) \quad A^{(1)} = (4, m)$$

$$z^{(2)} = (1, m) \quad A^{(2)} = (1, m)$$

$$d\hat{z}^{(1)} = (4, m) \quad dA^{(1)} = (4, m)$$

$$d\hat{z}^{(2)} = (1, m) \quad dA^{(2)} = (1, m)$$

dec 10: computing gradients:-

logistic regression:

$$\begin{array}{c}
 \text{Diagram showing the computation graph:} \\
 \begin{array}{c}
 \text{x} \rightarrow z = w^T x + b \\
 \text{d}w \quad w \quad \text{d}z \\
 \text{d}b \quad b
 \end{array}
 \xrightarrow{\quad a = \sigma(z) \quad}
 \boxed{a = \sigma(z)} \xrightarrow{\quad d(a|y) \quad}
 \boxed{d(a|y)}
 \end{array}$$

$$d(a|y) = -y \log a - (1-y) \log(1-a)$$

$$da = -\frac{y}{a} + \frac{(1-y)}{1-a}$$

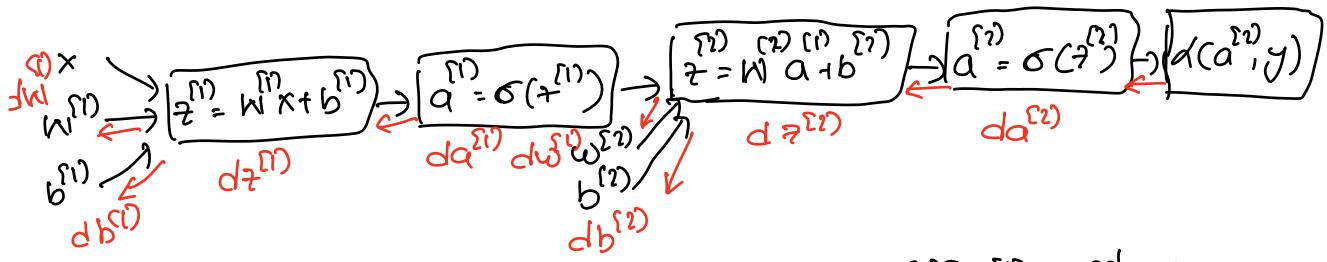
$$\begin{aligned}
 dz &= \frac{dL}{dz} = \frac{dL}{da} \frac{da}{dz} \\
 &= da \frac{d\sigma(z)}{dz} = da a(1-a)
 \end{aligned}$$

$$dz = \left(-\frac{y}{a} + \frac{(1-y)}{1-a}\right) a(1-a) = -y(1-a) + (1-y)a$$

$$dz = -y + ya + a - ya = a - y$$

$$\boxed{dz = a - y}$$

Gradient for one hidden layer - backprop



$$d\hat{z}^{(2)} = a^{(2)} - y$$

$$d\hat{z}^{(1)} = W^{(1)T} d\hat{z}^{(2)} * g'(z^{(1)})$$

$$dw^{(2)} = d\hat{z}^{(2)} a^{(1)T}$$

$$db^{(2)} = d\hat{z}^{(2)}$$

$$W^{(1)} = (n^{(1)}, n^{(1)})$$

$$z^{(1)}, d\hat{z}^{(1)} = (n^{(1)}, 1) = (1, 1)$$

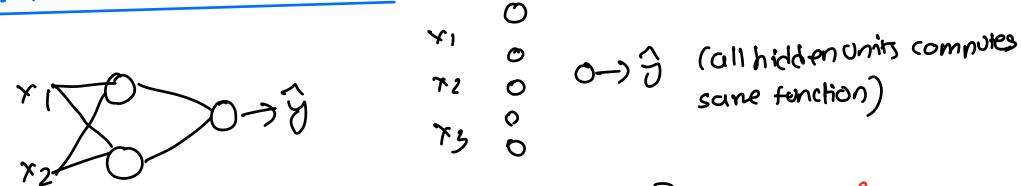
$$\hat{z}^{(1)}, d\hat{z}^{(1)} = (n^{(1)}, 1)$$

$$\begin{array}{l} x_1 \\ x_2 \\ x_3 \\ \vdots \\ n_x = n \end{array} \quad \begin{array}{l} 0 \\ 0 \\ 0 \\ \vdots \\ n^{(1)} = 1 \end{array} \quad 0 \rightarrow \hat{y}$$

$$d\hat{z}^{(1)} = W^{(2)T} d\hat{z}^{(2)} * g'(z^{(1)})$$

$$(n^{(1)}, 1) \quad (n^{(1)}, n^{(1)}) (n^{(1)}, 1) * (n^{(1)}, 1) = (n^{(1)}, 1)$$

dec12: Random Initialization



$$n^{(0)} = 2 \quad n^{(1)} = 2$$

$$W^{(1)} = (2, 2) = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \quad b^{(1)} = (2, 1) = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

problem?

$$W^{(2)} = [0, 0]$$

$$a_1^{(1)} = a_2^{(1)} = \text{same}$$

$$d\tilde{z}_1^{(1)} = d\tilde{z}_2^{(1)} = \text{same}$$

$$dW = \begin{bmatrix} v & v \\ v & v \end{bmatrix} \quad W^{(1)} = W^{(1)} - \alpha dW^{(1)}$$

if $\alpha^{(1)}$ is big

Grad will be slow

sol:

$$W^{(1)} = \text{np.random.randn}(2, 2) * 0.01$$

$$b^{(1)} = \text{np.zeros}(2, 1)$$

$$W^{(2)} = \text{np.random.randn}(1, 2) * 0.01$$

$$b^{(2)} = \text{np.zeros}(1, 1)$$

