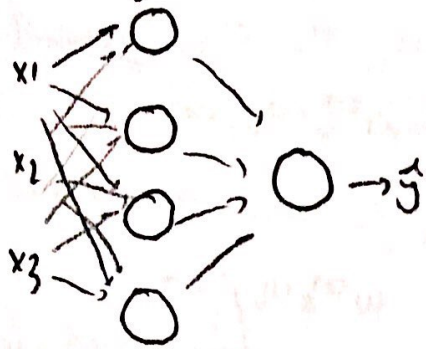


Vectorizing across multiple examples



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

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

$$z^{(2)} = w^{(2)}a^{(1)} + b^{(2)}$$

$$a^{(2)} = \sigma(z^{(2)})$$

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

$a^{(2)(i)}$ → example i
layer 2

→ for $i=1$ to n ,

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

$$a^{(1)(i)} = \sigma(z^{(1)(i)})$$

$$z^{(2)(i)} = w^{(2)}a^{(1)(i)} + b^{(2)}$$

$$a^{(2)(i)} = \sigma(z^{(2)(i)})$$

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

$(n \times m)$

$$z^{(1)} = w^{(1)}X + b^{(1)}$$

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

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

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

training examples
hidden units

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

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

ication for vectorized implementation

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

$$z^{(1)} x^{(2)} = w^{(1)} x^{(2)} + b^{(1)}$$

$$z^{(1)} x^{(3)} = w^{(1)} x^{(3)} + b^{(1)}$$

$$w^{(1)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

$$w^{(1)} x^{(1)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

$$w^{(1)} x^{(2)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

$$w^{(1)} x^{(3)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

$$w^{(1)} \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & x^{(3)} \\ | & | & | \end{bmatrix} = \begin{bmatrix} | & | & | \\ z^{(1)} x^{(1)} & z^{(1)} x^{(2)} & z^{(1)} x^{(3)} \\ | & | & | \end{bmatrix} = z^{(1)}$$

Gradient descent for neural networks

Parameters: $W^{(1)}, b^{(1)}, W^{(2)}, b^{(2)}$
 $(n^{(1)}, n^{(2)})$ $(n^{(1)}, 1)$ $(n^{(2)}, n^{(1)})$ $(n^{(2)}, 1)$

$$n_x = n^{(1)}, n^{(1)}, n^{(2)} = 1$$

Cost function: $J(W^{(1)}, b^{(1)}, W^{(2)}, b^{(2)}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(y, y)$
 \uparrow
 $a^{(2)}$

Gradient descent:

→ repeat { compute pred. $(y^{(1)}, i=1, \dots, n)$

$$dW^{(1)} = \frac{\partial J}{\partial W^{(1)}}, \quad db^{(1)} = \frac{\partial J}{\partial b^{(1)}}$$

$$W^{(1)} = W^{(1)} - \alpha dW^{(1)}$$

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



Formulas for computing derivatives

Forward propagation:

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

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

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

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

Backward propagation

$$dz^{(1)} = A^{(2)} - y$$

$$dW^{(2)} = \frac{1}{n} dz^{(2)} A^{(1)T}$$

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

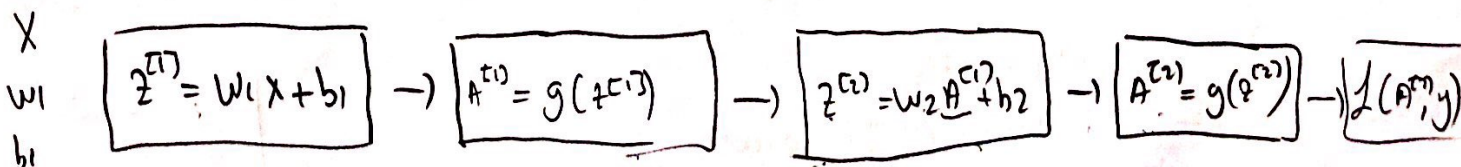
$$dz^{(1)} = W^{(2)T} dz^{(2)} \times g^{(1)}(z^{(1)})$$

elementwise prod $(n^{(1)}, n)$

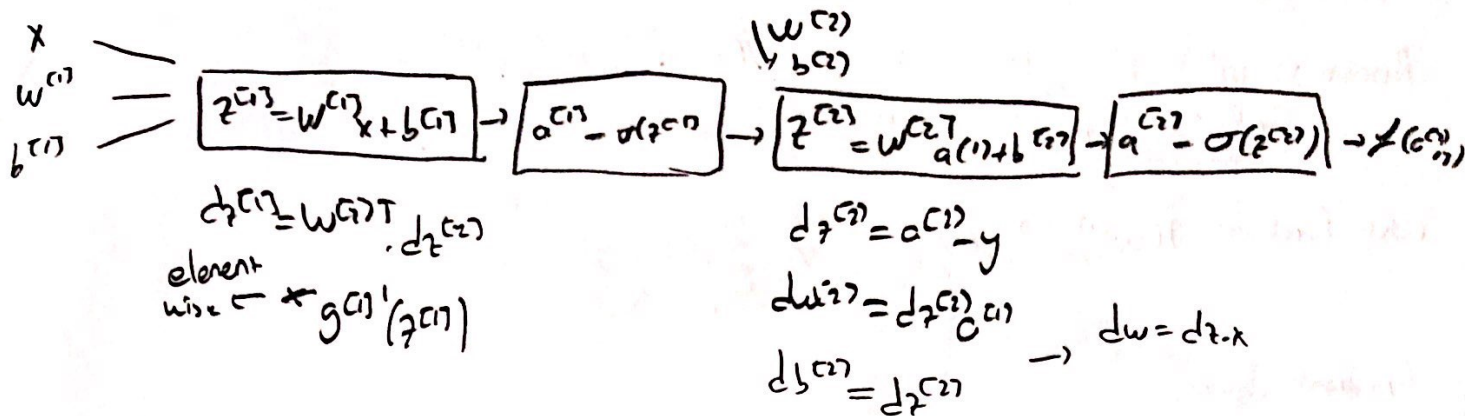
$$dW^{(1)} = \frac{1}{n} dz^{(1)} x^T$$

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

$$\begin{aligned} dz &= a - y \\ dw &= \frac{\partial J}{\partial w} \\ db &= \frac{\partial J}{\partial b} \end{aligned}$$



Set up gradients



$\rightarrow z^{(1)}, \delta z^{(1)}$

$z^{(1)} \delta z^{(1)}$

$$\delta w^{(1)} = \delta z^{(1)} x^T$$

$$\delta b^{(1)} = \delta z^{(1)}$$

$$\delta z^{(1)} = w^{(1)T} \delta z^{(2)} \odot y^{(1)'}(z^{(1)})$$

$$(n^{(1)}, 1) \pm n^{(1)}, n^{(1)}, (n^{(1)}, 1) \odot n^{(1)}, 1)$$

Summary of gradient descent

$$\delta z^{(1)} = a^{(1)} - y$$

$$\delta w^{(1)} = \delta z^{(1)} a^{(1)T}$$

$$\delta b^{(1)} = \delta z^{(1)}$$

$$\delta z^{(1)} = w^{(1)T} \delta z^{(2)} \odot y^{(1)'}(z^{(1)})$$

$$\delta w^{(1)} = \delta z^{(1)} x^T$$

$$\delta b^{(1)} = \delta z^{(1)}$$

$$\delta z^{(1)} = A^{(1)} - y$$

$$\delta w^{(1)} = \frac{1}{n} \delta z^{(1)} A^{(1)T}$$

$$\delta b^{(1)} = \frac{1}{n} \text{np.sum}(\delta z^{(1)}, \text{axis}=1, \text{keepdims}=True)$$

$$\delta z^{(1)} = w^{(1)T} \delta z^{(2)} \odot g^{(1)'}(z^{(1)})$$

$$\delta w^{(1)} = \frac{1}{n} \delta z^{(1)} x^T$$

$$\delta b^{(1)} = \frac{1}{n} \text{np.sum}(\delta z^{(1)}, \text{axis}=1, \text{keepdims}=True)$$

Random Initialization

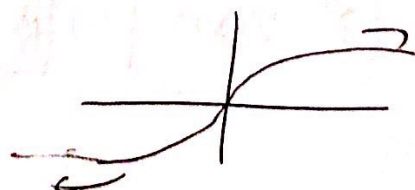
$w^{(1)} = \text{np.random.randn}(2, 1) \times 0.001$ \rightarrow gaussian

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

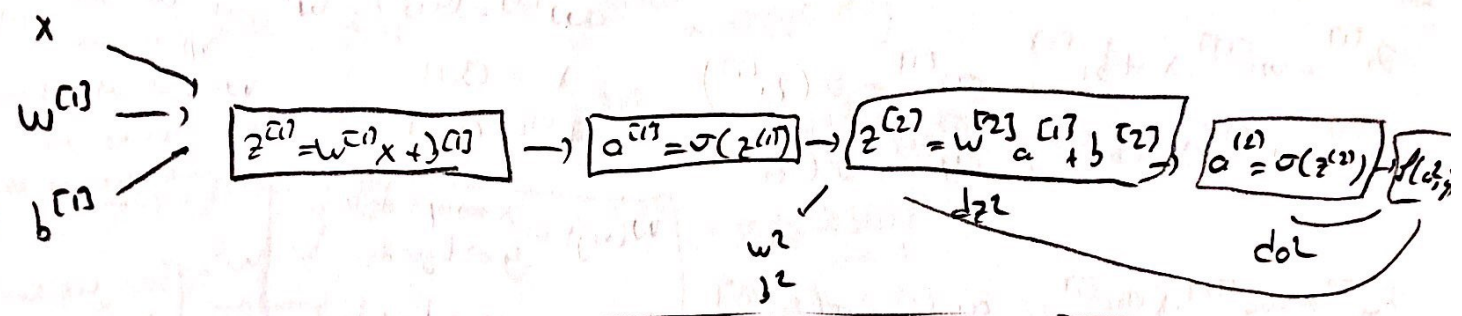
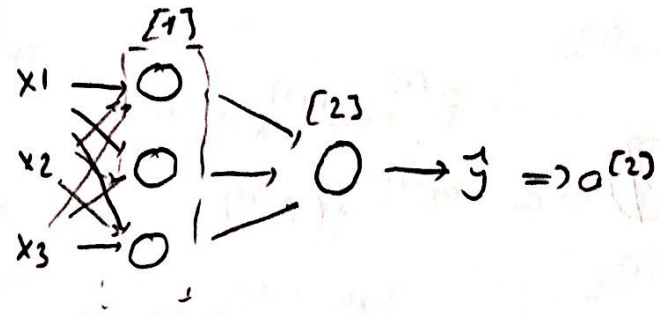
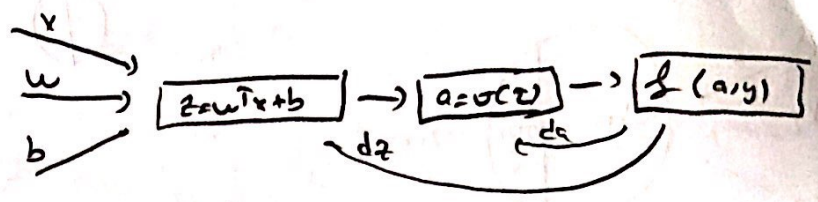
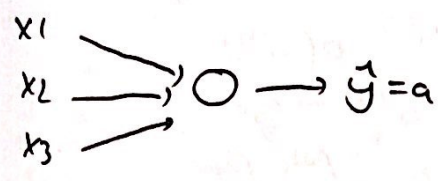
$w^{(1)} = \text{np.random.randn}(1, 1) \times 0.001$

$b^{(1)} = 0$

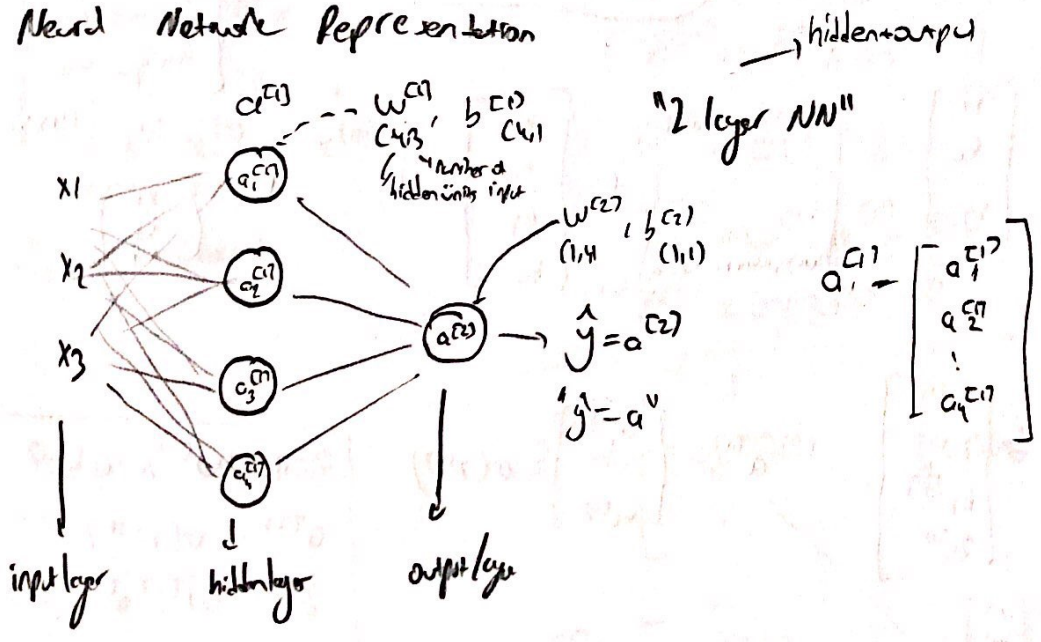
Make ten aggressive scale down
the bias we could have degree 1 way
-1 e solution



What is a Neural Network?



Neural Network Representation



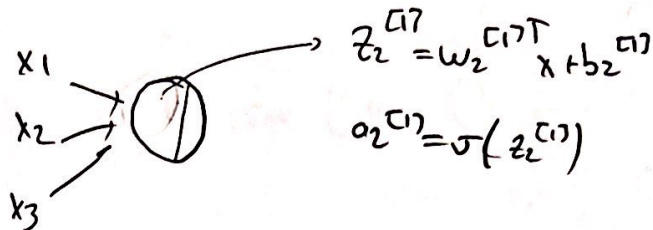
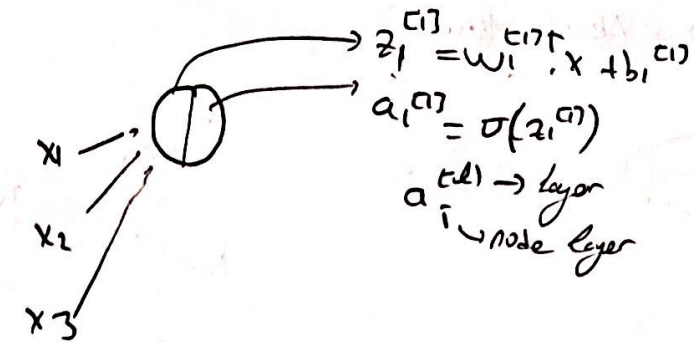
$$a^{(1)} = \begin{bmatrix} a_1^{(1)} \\ a_2^{(1)} \\ a_3^{(1)} \\ a_4^{(1)} \end{bmatrix}$$

Network Representation



$$z = w^T x + b$$

$$a = \sigma(z)$$



$$z_1^{(l)} = w_1^{(l)T} \cdot x + b_1^{(l)}, \quad a_1^{(l)} = \sigma(z_1^{(l)})$$

$$z_2^{(l)} = w_2^{(l)T} \cdot x + b_2^{(l)}, \quad a_2^{(l)} = \sigma(z_2^{(l)})$$

...

$$z_n^{(l)} = w_n^{(l)T} \cdot x + b_n^{(l)}, \quad a_n^{(l)} = \sigma(z_n^{(l)})$$

$$w^{(l)} = (w_i) = w^T$$

$$x = (3, 1)$$

$$b^{(l)} = (4, 1)$$

$$w = \begin{bmatrix} w_1 & w_2 & w_3 & w_4 \\ w_1 & w_2 & w_3 & w_4 \\ w_1 & w_2 & w_3 & w_4 \\ w_1 & w_2 & w_3 & w_4 \end{bmatrix}_{4 \times 4}$$

$$w^T = \begin{bmatrix} w_1 & w_1 & w_1 & w_1 \\ w_2 & w_2 & w_2 & w_2 \\ w_3 & w_3 & w_3 & w_3 \\ w_4 & w_4 & w_4 & w_4 \end{bmatrix}_{(4 \times 4)}$$

$W(x, y) = x \Rightarrow$ input layer
 $y \Rightarrow$ layer deep
 $b(z, 1) \Rightarrow z \Rightarrow$ layer deep

$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$w^{(l)T} \cdot x = \begin{bmatrix} w_1 & w_1 & w_1 \\ w_2 & w_2 & w_2 \\ w_3 & w_3 & w_3 \\ w_4 & w_4 & w_4 \end{bmatrix}_{4 \times 3} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}_{3 \times 1} = \begin{bmatrix} w_1 x_1 + w_1 x_2 + w_1 x_3 \\ \vdots \\ w_4 x_1 + w_4 x_2 + w_4 x_3 \end{bmatrix}_{4 \times 1} + \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ b_4 \end{bmatrix}_{4 \times 1}$$

Transpose of w is w^T

$$= \begin{bmatrix} w_1^{(l)T} x + b_1^{(l)} \\ w_2^{(l)T} x + b_2^{(l)} \\ w_3^{(l)T} x + b_3^{(l)} \\ w_4^{(l)T} x + b_4^{(l)} \end{bmatrix} = \begin{bmatrix} z_1^{(l)} \\ z_2^{(l)} \\ z_3^{(l)} \\ z_4^{(l)} \end{bmatrix}$$

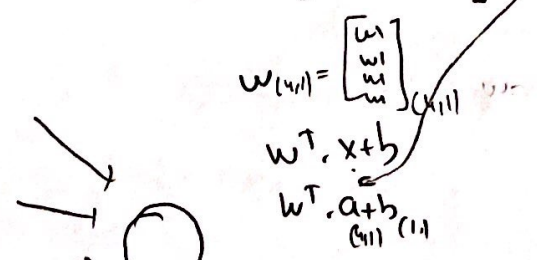
$$a^{(l)} = \begin{bmatrix} \sigma(z_1^{(l)}) \\ \sigma(z_2^{(l)}) \\ \sigma(z_3^{(l)}) \\ \sigma(z_4^{(l)}) \end{bmatrix} = \sigma(z^{(l)})$$

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

$$a^{(l)} = \sigma(z^{(l)})$$

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

$$a^{(l+1)} = \sigma(z^{(l+1)})$$



$$\begin{bmatrix} w_1 & w_1 & w_1 & w_1 \\ w_2 & w_2 & w_2 & w_2 \\ w_3 & w_3 & w_3 & w_3 \\ w_4 & w_4 & w_4 & w_4 \end{bmatrix} \begin{bmatrix} a_1^{(l)} \\ a_2^{(l)} \\ a_3^{(l)} \\ a_4^{(l)} \end{bmatrix} = (a_1^{(l)} \cdot w_1 + a_2^{(l)} \cdot w_2 + a_3^{(l)} \cdot w_3 + a_4^{(l)} \cdot w_4) + b = z^{(l+1)}$$

$$a^{(l+1)} = \sigma(z^{(l+1)})$$