

## Binary Classification

Cat (1) vs Non-Cat (0)

$$x = \begin{bmatrix} 255 \\ 231 \\ \vdots \end{bmatrix} \quad \text{Image} - \underbrace{\begin{bmatrix} 64 \times 64 \times 3 \end{bmatrix}^T}_{\text{Image dim}} = 12288 \quad n = n_x = 12288$$

Notat<sup>n</sup>:  $(x, y) \quad x \in \mathbb{R}^{n_x}, y \in \{0, 1\}$

$m$  training examples:  $\{(x^1, y^1), (x^2, y^2), \dots, (x^m, y^m)\}$

$$M = M_{\text{train}} \quad M_{\text{test}} = \# \text{test examples}$$

$$x = \begin{bmatrix} | & | & | \\ x^1 & x^2 & \cdots & x^m \\ | & | & \cdots & | \end{bmatrix}^T \quad x \in \mathbb{R}^{n_x \times m} \quad \left[ \begin{array}{l} \text{X. shape} = (n_x, m) \\ \text{Y. shape} = (1, m) \end{array} \right]$$

$$y = [y^1, y^2, \dots, y^m] \quad y \in \mathbb{R}^{1 \times m} \quad \left[ \begin{array}{l} \text{Y. shape} = (1, m) \end{array} \right]$$

## LOGISTIC REGRESSION

Given  $x$ , want  $\hat{y} = P(y=1|x)$  where  $0 \leq \hat{y} \leq 1$

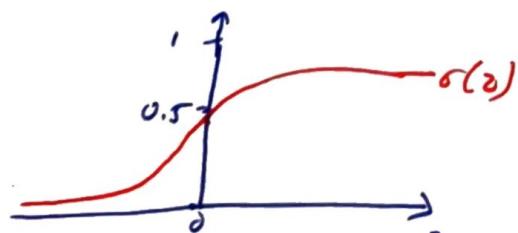
Parameters:  $w \in \mathbb{R}^{n_x}$ ,  $b \in \mathbb{R}$

Output  $\hat{y} = w^T x + b$  }  $\left\{ \begin{array}{l} \text{it would be difficult to predict :: here} \\ \hat{y} \text{ might be } > 1 \text{ or } < 0 \end{array} \right.$

$$\therefore \hat{y} = \sigma(w^T x + b)$$

↓  
sigmoid

$$\boxed{\sigma(z) = \frac{1}{1 + e^{-z}}}$$



If  $z$  large  $\Rightarrow \sigma(z) \approx \frac{1}{1+0} = 1$

If  $z$  large -ve no.  $\Rightarrow \sigma(z) = \frac{1}{1+\text{big no.}} \approx 0$

$$\hat{y}^i = \sigma(\omega^T x^i + b), \text{ where } \sigma(z^i) = \frac{1}{1+e^{-z^i}} \quad z^i = \omega^T x^i + b$$

Given  $\{(x^1, y^1), \dots, (x^m, y^m)\}$ , want  $\hat{y}^i \approx y^i$

Loss (error) f:

$$L(\hat{y}, y) = - (y \log \hat{y} + (1-y) \log (1-\hat{y}))$$

If  $y=1$ , then  $L(\hat{y}, y) = -\log \hat{y}$   $\begin{cases} \Rightarrow \text{want } \log \hat{y} \text{ to be large,} \\ \therefore \text{want } \hat{y} \text{ large} \\ \text{but since it is binary classification problem, } \hat{y} \text{ cannot be } > 1 \\ \Rightarrow \text{we want } \hat{y} \text{ to be as close to } 1 \end{cases}$

If  $y=0$ , then  $L(\hat{y}, y) = -\log(1-\hat{y}) \Rightarrow \begin{cases} \text{want } \log(1-\hat{y}) \text{ to be large,} \\ \text{only want } \hat{y} \text{ small} \\ \text{only want } \hat{y} \text{ to be as close to } 0 \end{cases}$

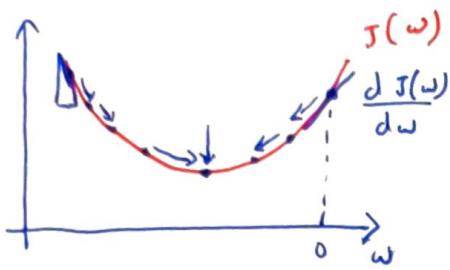
Cost f:

$$J(\omega, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^i, y^i) = \frac{1}{m} \sum_{i=1}^m [y^i \log \hat{y}^i + (1-y^i) \log(1-\hat{y}^i)]$$

or  $J(\omega, b) = -\frac{1}{m} \sum_{i=1}^m [y^i \log \hat{y}^i + (1-y^i) \log(1-\hat{y}^i)]$

## Gradient Descent:

Want to find  $w, b$  that minimize  $J(w, b)$   $\rightarrow$  convex f



Repeat: {  
 $w: w - \alpha \frac{dJ(w)}{dw}$  learning rate  
 }

~ly:  $J(w, b)$

Repeat {

$$w: w - \alpha \frac{\partial J(w, b)}{\partial w}$$

$$b: b - \alpha \frac{\partial J(w, b)}{\partial b}$$

}

## Computation Graph

$$J(a, b, c) = 3(a + bc)$$

$$u = bc$$

$$v = a + u$$

$$J = 3v$$

$$\frac{dJ}{da} = ?$$

$$a = 5 \rightarrow 5.001$$

$$v = a + u = 11 \rightarrow 11.001$$

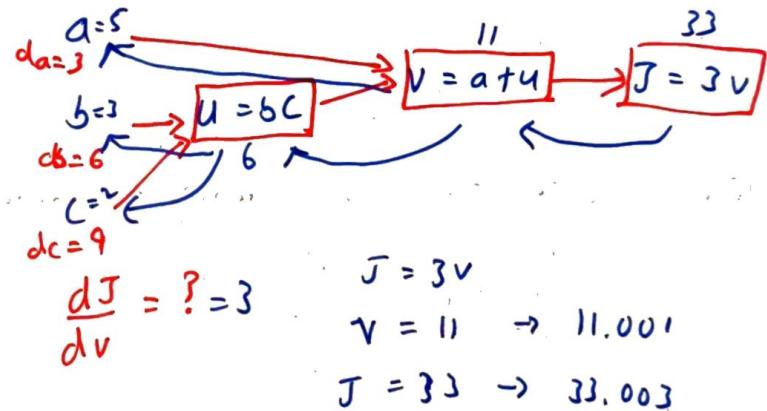
$$J = 33 \rightarrow 33.003$$

$$\Rightarrow \frac{dJ}{da} = 3 = \frac{dJ}{dv} \cdot \frac{dv}{da}$$

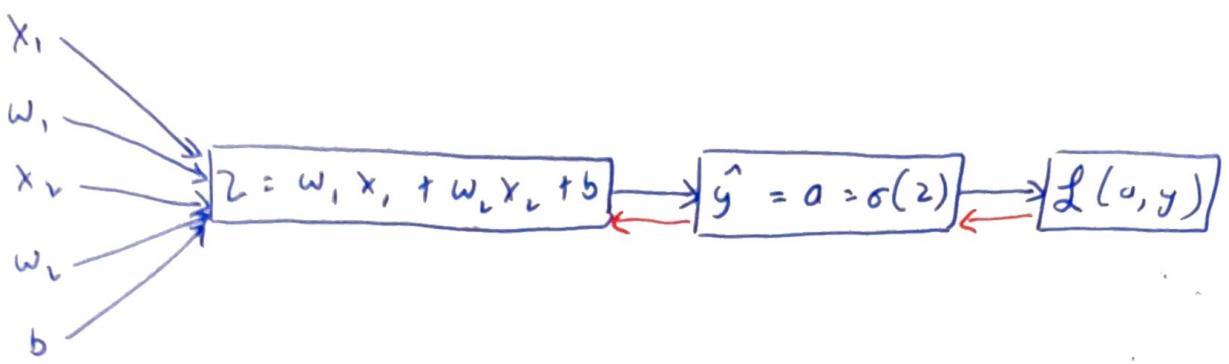
$$\frac{dv}{da} = 1 \quad \text{"do"}$$

$$\text{~ly } \frac{dJ}{du} = 3 \{ \text{"du"} \}, \quad \frac{dJ}{db} = \frac{dJ}{du} \frac{du}{db} = 3 \times 2 = 6 \{ \text{"db"} \}$$

$$\text{~ly "dc" = 9"}$$



$\frac{d(\text{Final Output Variable})}{d(\text{var})} \Rightarrow \text{in code}$   
 $d(\text{var})$  "d var"



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

$$\begin{aligned} d_z &= \frac{dL}{dz} = \frac{dL}{da} \cdot \frac{da}{dz} = \left(-\frac{y}{a} + \frac{1-y}{1-a}\right) \cdot (a(1-a)) \\ &= a(y - a) \end{aligned}$$

$$\begin{aligned} d_{w_1} &= \frac{\partial L}{\partial w_1} = x_1 \cdot d_z \Rightarrow \\ &\quad \begin{cases} w_1 := w_1 - \alpha d_{w_1} \\ w_2 := w_2 - \alpha d_{w_2} \\ b := b - \alpha d_b \end{cases} \\ d_{w_2} &= x_2 \cdot d_z \\ d_b &= d_z \end{aligned}$$

Gradient Descent on m training examples:

$$\begin{aligned} J(w, b) &= \frac{1}{m} \sum_{i=1}^m L(a^i, b^i) \\ \rightarrow a^i &= \hat{y}^i = \sigma(z^i) = \sigma(w^T x^i + b) \\ \frac{\partial}{\partial w_i} J(w, b) &= \underbrace{\frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial w_i} L(a^i, y^i)}_{d_{w_i}^i} = (x^i, y^i) \end{aligned}$$

Let,  $J=0, d_{w_1}=0, d_{w_2}=0, d_b=0$

For  $i = 1$  to  $m$

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

$$a^i = \sigma(z^i)$$

$$J_t = -[y^i \log a^i + (1-y^i) \log (1-a^i)]$$

$$d_z^i = a^i - y^i$$

$$\begin{aligned} \frac{d w_i}{d z^i} &= x_i \\ \frac{d w_i}{d z^i} &= x_i^i d_{z^i} \quad \uparrow n=2 \\ \frac{d b}{d z^i} &= d_{z^i} \end{aligned}$$

$$J/m = m$$

$$d\omega_1/m ; d\omega_2/m ; db/m$$

$$\Rightarrow d\omega_i = \frac{\partial J}{\partial \omega_i}$$

$$\omega_1 := \omega_1 - \alpha d\omega_1$$

$$\omega_2 := \omega_2 - \alpha d\omega_2$$

$$b := b - \alpha db$$

Issue: :: L for loops for m training sets  
for each feature \Rightarrow making computation complex

$\Rightarrow$  Solution

### VECTORIZATION

$$z = \omega^T x + b$$

$$\omega = \begin{bmatrix} \cdot \\ \vdots \\ \cdot \end{bmatrix} \quad x = \begin{bmatrix} \cdot \\ \vdots \\ \cdot \end{bmatrix} \quad \begin{array}{l} \omega \in \mathbb{R}^{n_x} \\ x \in \mathbb{R}^{n_x} \end{array}$$

Non Vectorized:

$$\begin{aligned} z &= 0 \\ \text{for } i \text{ in range } (n-x): \\ z_t &= \omega[i] * x(i) \\ z_t &= b \end{aligned}$$

Vectorized:

$$z = \underbrace{\text{np. dot}(\omega, x)}_{\omega^T x} + b$$

Logistic regression derivatives (vectorized)

$$J = 0, d\omega = np.zeros(n_x, 1), db = 0$$

for  $i = 1$  to  $m$ :

$$z^i = \omega^T x^i + b$$

$$a^i = \sigma(z^i)$$

$$J += -[y^i \log a^i + (1-y)(\log(1-a^i))]$$

$$d\omega^i = a^i(1-a^i)$$

$$d\omega += x^i d\omega^i$$

$$db += d\omega^i$$

$$J = J/m, d\omega / \alpha = m, db = db/m$$

## VECTORIZING LOGISTIC REGRESSION

$$z^{(1)} = \omega^T x^{(1)} + b$$

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

$$z^{(2)} = \omega^T x^{(2)} + b$$

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

$$z^{(3)} = \omega^T x^{(3)} + b$$

$$a^{(3)} = \sigma(z^{(3)})$$

$$X = \begin{bmatrix} 1 & x^{(1)} & 1 \\ 1 & x^{(2)} & 1 \\ \vdots & \vdots & \vdots \\ 1 & x^{(m)} & 1 \end{bmatrix} \in \mathbb{R}^{n_x \times m}$$

$$\begin{bmatrix} z^{(1)} & z^{(2)} & \dots & z^{(m)} \end{bmatrix} = \omega^T X_T \begin{bmatrix} b & b & \dots & b \end{bmatrix}_{1 \times m}$$

$$= \underbrace{\omega^T \begin{bmatrix} 1 & x^{(1)} & 1 \\ 1 & x^{(2)} & 1 \\ \vdots & \vdots & \vdots \\ 1 & x^{(m)} & 1 \end{bmatrix}}_{m \times n} \begin{bmatrix} b & b & \dots & b \end{bmatrix}_{1 \times m}$$

$$Z = \begin{bmatrix} \omega^T x^{(1)} + b & \omega^T x^{(2)} + b & \dots & \omega^T x^{(m)} + b \end{bmatrix}$$

<sup>Algorithm</sup>  
 $\{Z = np.\text{dot}(\omega.T, x) + b\} \Rightarrow \text{"Broadcasting"}$

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

## VECTORIZING LOGISTIC REGRESSION'S GRADIENT DESCENT

$$d_2^{(1)} = a^{(1)} - y^{(1)}$$

$$d_2^{(2)} = a^{(2)} - y^{(2)}$$

$$dZ = \begin{bmatrix} d_2^{(1)} & d_2^{(2)} & \dots & d_2^{(m)} \end{bmatrix}_{1 \times m}$$

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

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

$$dw = 0$$

$$dw += x^{(1)} d_2^{(1)}$$

$$dw += x^{(2)} d_2^{(2)}$$

:

$$dw / = M$$

$$db = 0$$

$$db += d_2^{(1)}$$

$$db += d_2^{(2)}$$

:

$$db += d_2^{(m)}$$

$$db / = m$$

$$\Rightarrow \delta b = \frac{1}{m} \sum_{i=1}^m d_2^{(i)}$$

$$\underline{\delta b = \frac{1}{m} \text{np.sum}(d_2)}$$

nly

$$\begin{aligned} \frac{dw}{\delta w} &= \frac{1}{m} X d_2^T \\ &= \frac{1}{m} \begin{bmatrix} 1 & 1 & \dots & 1 \end{bmatrix} \begin{bmatrix} d_2^{(1)} \\ d_2^{(2)} \\ \vdots \\ d_2^{(m)} \end{bmatrix} \\ &= \frac{1}{m} [x^{(1)} d_2^{(1)} + \dots + x^{(m)} d_2^{(m)}] \end{aligned}$$

## Implementing Logistic Regression

$$Z = \omega^T X + b$$

$$A = \sigma(Z)$$

$$= \text{np.dot}(\omega, X) + b$$

$$d_2 = A - Y$$

$$dw = \frac{1}{m} X d_2^T$$

$$\delta b = \frac{1}{m} \text{np.sum}(d_2)$$

$$\omega = \omega - \alpha dw$$

$$b = b - \alpha \delta b$$

## 18 road Carting Example

Calories from Carbs, Proteins, Fats in 100g of diff. foods:

	Apples	Beef	Eggs	Potatoes	
Carb	56.0	0.0	4.4	68.0	
Protein	1.2	104.0	52.0	8.0	
Fat	1.8	135.0	99.0	0.9	

$A_{3 \times 4}$

$(3 \times 4)$

Calculate % of Calories from Carbs, Potatoes, Fat. Can you do this without explicit for-loop?

$$A = \text{np.array}([[ \dots ], [\dots], [\dots], [\dots]])$$

$$col = A.sum(axis=0)$$

~~$$\text{percentage} = A / col \cdot \text{percentage}(+, \%) * 100$$~~

$$\text{percentage} = A / col \cdot 100$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}_{m \times n} + \begin{bmatrix} 100 & 200 & 300 \end{bmatrix}_{1 \times n} = \begin{bmatrix} 101 & 202 & 303 \\ 104 & 205 & 306 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}_{m \times n} + \begin{bmatrix} 100 \\ 200 \end{bmatrix} = \begin{bmatrix} 101 & 102 & 103 \\ 204 & 205 & 206 \end{bmatrix}_{m \times n}$$

## General Principle:

$$\begin{matrix} (m, n) \\ \text{matrix} \end{matrix} \quad \begin{matrix} + \\ * \\ / \end{matrix} \quad \begin{matrix} (1, n) \rightsquigarrow (m, n) \\ (m, 1) \rightsquigarrow (m, n) \end{matrix}$$

$$\begin{matrix} (m, 1) \\ (1, n) \end{matrix} \quad \begin{matrix} + \\ * \\ / \end{matrix} \quad \begin{matrix} R_{(1,1)} \rightsquigarrow (m, 1) \\ R_{(1,1)} \rightsquigarrow (1, n) \end{matrix}$$

## Note:

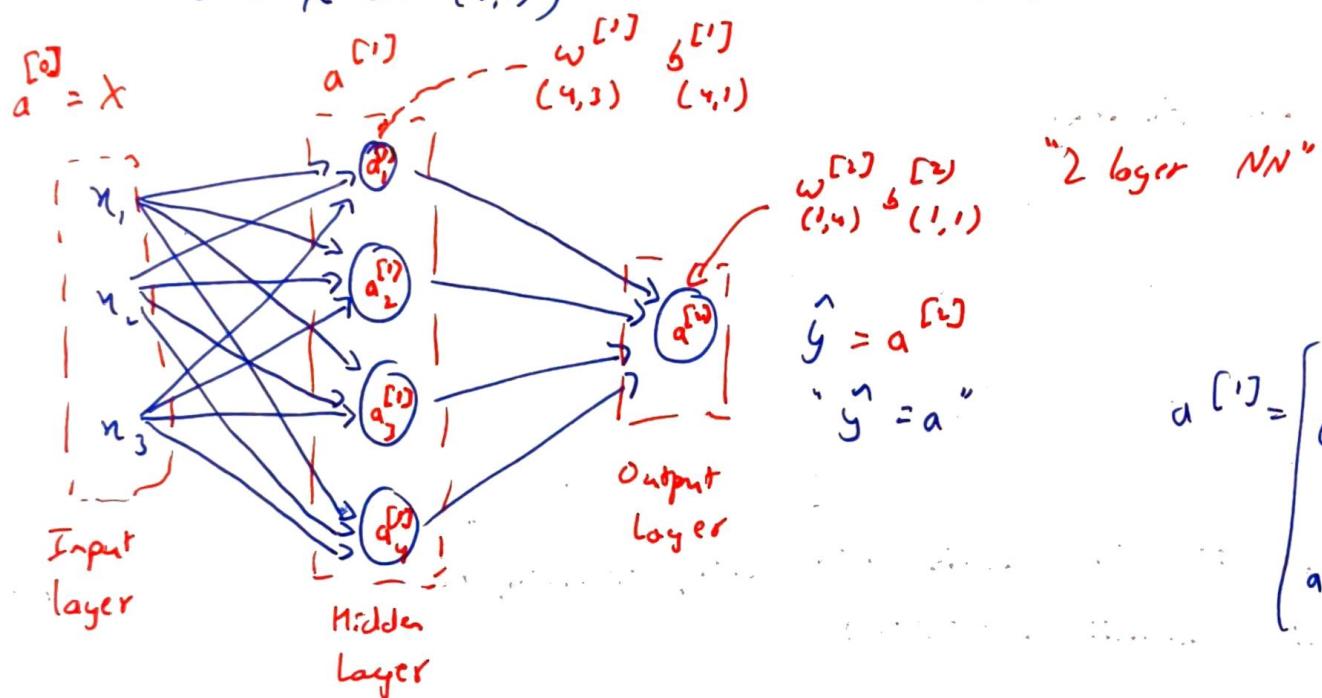
i)  $a = np.random.rand(5)$

$a.shape : (5,) \Rightarrow \text{"rank 1 array"}$  } - Don't use

$a = np.random.rand(5, 1) \rightarrow \text{shape } (5, 1) \Rightarrow \text{column vector}$

$a = np.random.rand(1, 5) \rightarrow \text{shape } (1, 5) \Rightarrow \text{row vector}$

assert (a.shape == (5, 1))



$$a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ \vdots \\ a_4^{[1]} \end{bmatrix}$$

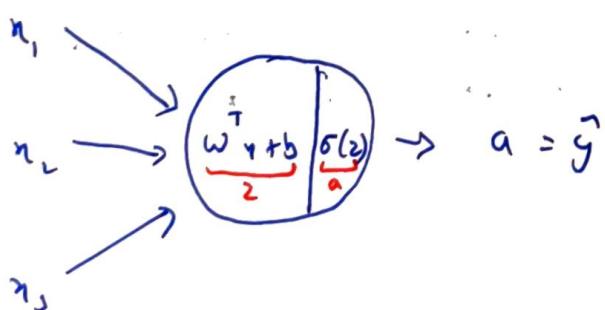
$$z_1^{[1]} = w_1^{[1]T} + b_1^{[1]}$$

$$a_1^{[1]} = \sigma(z_1^{[1]})$$

$a_i^{[1]}$  → layer

$a_i^{[1]}$  → node in layer

$$z_2^{[1]} = w_2^{[1]T} + b_2^{[1]}$$



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

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

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

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

$$\boxed{Z = \begin{bmatrix} \omega_1^{(1)T} \\ \omega_2^{(1)T} \\ \omega_3^{(1)T} \\ \omega_4^{(1)T} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} + \begin{bmatrix} b_1^{(1)} \\ b_2^{(1)} \\ b_3^{(1)} \\ b_4^{(1)} \end{bmatrix}} = \begin{bmatrix} \omega_1^{(1)T} x + b_1^{(1)} \\ \omega_2^{(1)T} x + b_2^{(1)} \\ \omega_3^{(1)T} x + b_3^{(1)} \\ \omega_4^{(1)T} x + b_4^{(1)} \end{bmatrix} = \begin{bmatrix} z_1^{(1)} \\ z_2^{(1)} \\ z_3^{(1)} \\ z_4^{(1)} \end{bmatrix}$$

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

Given Input  $x$  ( $a^{(0)}$ ):

$$\boxed{\begin{aligned} \rightarrow z^{(1)} &= \omega^{(1)T} a^{(0)} + b^{(1)} \\ \rightarrow a^{(1)} &= \sigma(z^{(1)}) \\ \rightarrow z_{1x1}^{(2)} &= \omega_{1x1}^{(2)} a_{1x1}^{(1)} + b_{1x1}^{(2)} \\ \rightarrow a_{1x1}^{(2)} &= \sigma(z^{(2)})_{1x1} \end{aligned}}$$

## Vectorizing across multiple examples:

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

$a^{[2](i)}$   $\downarrow$   $\rightarrow$  training eg. i  
layer 2

for  $i = 1$  to  $m$ :

$$\left\{ \begin{array}{l} z^{[1](i)} = w^{(1)} x^i + b^{(1)} \\ a^{[1](i)} = \sigma(z^{[1](i)}) \\ z^{[2](i)} = w^{[2]} a^{[1](i)} + b^{[2]} \\ a^{[2](i)} = \sigma(z^{[2](i)}) \end{array} \right.$$

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

To compute:

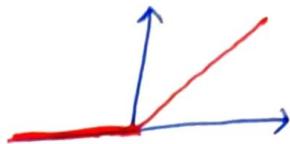
$$\begin{aligned} z^{[1]} &= w^{(1)} X + b^{(1)} \\ A^{[1]} &= \sigma(z^{[1]}) \\ z^{[2]} &= w^{[2]} A^{[1]} + b^{[2]} \\ A^{[2]} &= \sigma(z^{[2]}) \end{aligned}$$

$$z^{[1]} = \begin{bmatrix} | & | & | \\ z^{[1](1)} & z^{[1](2)} & \dots & z^{[1](m)} \\ | & | & & | \end{bmatrix}$$

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

hidden units.  $\xrightarrow{\text{Training egs.}}$

**ReLU** Rectified Linear Unit  
 $a = \max(0, z)$



Note:

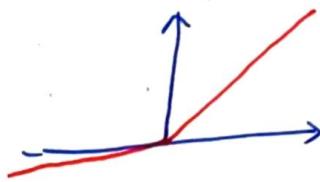
If the value of output is (0,1)  $\rightarrow$  sigmoid activation in output unit

ReLU - default

Advantage  $\rightarrow$  derivative = 0 when -ve value

Leaky ReLU:

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



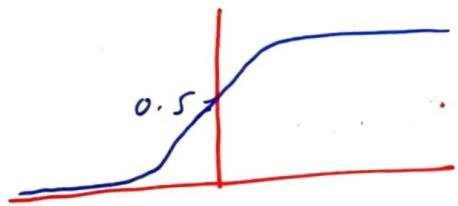
Derivations of Activation F:

i) Sigmoid activat<sup>n</sup> f<sup>n</sup>:

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

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

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



If  $z = 10$ ,  $g(z) \approx 1$

$$\frac{dg(z)}{dz} \approx 1(1-1) \approx 0$$

When  $z = -10$ ,  $g(z) \approx 0$

$$\frac{dg(z)}{dz} \approx 0(1-0) \approx 0$$

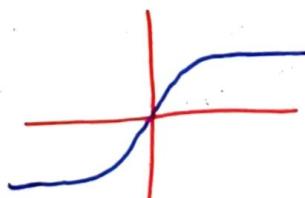
When  $z = 0$ ,  $g(z) = \frac{1}{2}$

$$\frac{dg(z)}{dz} = \frac{1}{2} \left( 1 - \frac{1}{2} \right) = \frac{1}{4}$$

v) Tanh activat<sup>n</sup> f<sup>n</sup>:

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

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



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

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

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

## 3) ReLU

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

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

## 4) Leaky ReLU

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

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

GRADIENT DESCENT FOR NEURAL NETWORKS

Parameters:  $\omega^{[1]}_{n^{[1]}, n^{[2]}}$ ,  $b^{[1]}_{n^{[1]}, 1}$ ,  $\omega^{[2]}_{n^{[2]}, n^{[1]}}$ ,  $b^{[2]}_{n^{[2]}, 1}$

$$h_x = a^{[0]}, a^{[1]}, a^{[2]} = 1$$

$$\text{Cost } f: J(\omega^{[1]}, b^{[1]}, \omega^{[2]}, b^{[2]}) = \frac{1}{m} \sum_{i=1}^m \hat{L}(y_i, \hat{y}_i)$$

Gradient descent:

Repeat {

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

$$d\omega^{[1]} = \frac{dJ}{d\omega^{[1]}}, db^{[1]} = \frac{dJ}{db^{[1]}}, \dots$$

$$\omega^{[1]} = \omega^{[1]} + \alpha d\omega^{[1]}$$

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

$$\vdots \quad ; \quad \vdots$$

Forward Propogation

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

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

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

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

Back propagation:

$$d_2^{[2]} = A^{[2]} - Y$$

$$d\omega^{[2]} = \frac{1}{m} d_2^{[2]} A^{[1]T}$$

$$db^{[2]} = \frac{1}{m} \text{np.sum}(d_2^{[2]}, \text{axis}=1, \text{keepdims=True})$$

parents from for  
list may like

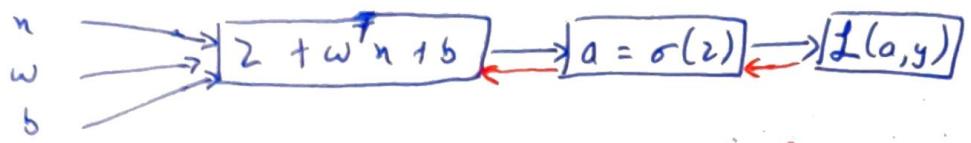
$$dz^{[1]} = \underbrace{\omega^{[1]T} d_2^{[2]}}_{(n^{[1]}, n^{[2]})} * \underbrace{g^{[1]}'(z^{[1]})}_{\text{element wise product}}$$

$$d\omega^{[1]} = \frac{1}{m} dz^{[1]} X^T$$

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

## Backpropagation Intuition

Logistic regression:



$$da = \frac{d}{da} L(a, y) = -y \log a - (1-y) \log(1-a)$$

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

$$dz = a - y$$

$$dz = da \cdot g'(z) \quad \left\{ \begin{matrix} \text{using } g(z) = \sigma(z) \\ dz = da \cdot \sigma'(z) \end{matrix} \right.$$

$$\frac{\partial L}{\partial z} = \frac{\partial L}{\partial a} \cdot \boxed{\frac{da}{dz}} \rightarrow \frac{d}{dz} g(z) = g'(z)$$

$$dw = dz \cdot x$$

$$db = dz$$

Neural Network gradients:

A diagram of a neural network with two layers. Inputs  $x$  pass through layer 1 with weights  $w^{[1]}$  and bias  $b^{[1]}$  to produce  $z^{[1]} = w^{[1]^T} n^{[1]} + b^{[1]}$ . The output of layer 1 is  $a^{[1]} = \sigma(z^{[1]})$ . Inputs  $a^{[1]}$  pass through layer 2 with weights  $w^{[2]}$  and bias  $b^{[2]}$  to produce  $z^{[2]} = w^{[2]^T} a^{[1]} + b^{[2]}$ . The output of layer 2 is  $a^{[2]} = \sigma(z^{[2]})$ . The total loss is  $L(a^{[2]}, y)$ .

$$\text{dimensions: } w^{[2]} \rightarrow (n^{[2]}, n^{[1]})$$

$$z^{[2]}, dz^{[2]} \rightarrow (n^{[2]}, 1) - (1, 1)$$

$$z^{[1]}, dz^{[1]} \rightarrow (n^{[1]}, 1)$$

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

$$d\omega^{[1]} = d_2^{[1]} \cdot x^T$$

$$d_b^{[1]} = d_2^{[1]}$$

Summary:

$$\left\{ \begin{array}{l} d_2^{[2]} = a^{[2]} - y \\ d\omega^{[2]} = d_2^{[2]} a^{[1]T} \\ d_b^{[2]} = d_2^{[2]} \\ d_2^{[1]} = \omega^{[2]T} d_2^{[2]} * g^{[1]'}(z^{[1]}) \\ d\omega^{[1]} = d_2^{[1]} x^T \\ d_b^{[1]} = d_2^{[1]} \end{array} \right.$$

vectorized implementation?

$$\begin{aligned} z^{[1]} &= \omega^{[1]} x + b^{[1]} \\ a^{[1]} &= g^{[1]}(z^{[1]}) \end{aligned}$$

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

$$\begin{aligned} z^{[1]} &= \omega^{[1]} x + b^{[1]} \\ A^{[1]} &= g^{[1]}(z^{[1]}) \end{aligned}$$

→ Vectorized equations:

$$d_2^{[2]} = A^{[2]} - y$$

$$d\omega^{[2]} = \frac{1}{m} d_2^{[2]} A^{[1]T}$$

$$d_b^{[2]} = \frac{1}{n} \text{np.sum}(d_2^{[2]}, \text{axis}=1, \text{keepdims=True})$$

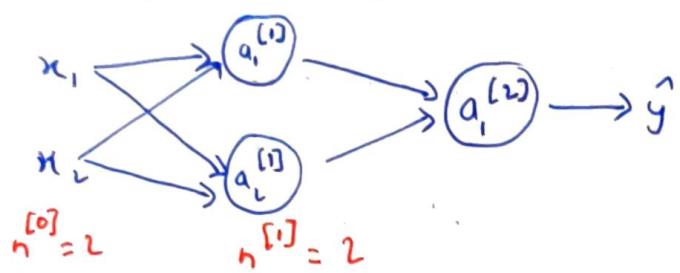
$$d_2^{[1]} = \underbrace{\omega^{[2]T} d_2^{[2]} * g^{[1]'}(z^{[1]})}_{(n^{[1]}, m) \text{ element wise product}} \quad (n^{[1]}, m)$$

$$d\omega^{[1]} = \frac{1}{m} d_2^{[1]} x^T$$

$$d_b^{[1]} = \frac{1}{m} \text{np.sum}(d_2^{[1]}, \text{axis}=1, \text{keepdims=True})$$

## Random Initialization:

What happens if you initialize weights to zeros?



$$\omega_{2 \times 2}^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}, \quad b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$\Rightarrow a_1^{[1]} = a_2^{[1]}$ , & when done back propagat', then  
 $d_2^{[1]} = d_2^{[1]}$

$\Rightarrow$  Symmetric neuron (Identical)

$$d\omega = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}, \quad \therefore \omega^{[1]} = \omega^{[1]} - \alpha d\omega$$

Then  $\because$  both  $u_k$  are same,  $d_2$  are same,  $\therefore$  weight will be updated in the same way  $\Rightarrow$  act as single same unit

What should be done?

$$\omega^{[1]} = \text{np.random.rand}(2,2) * 0.01$$

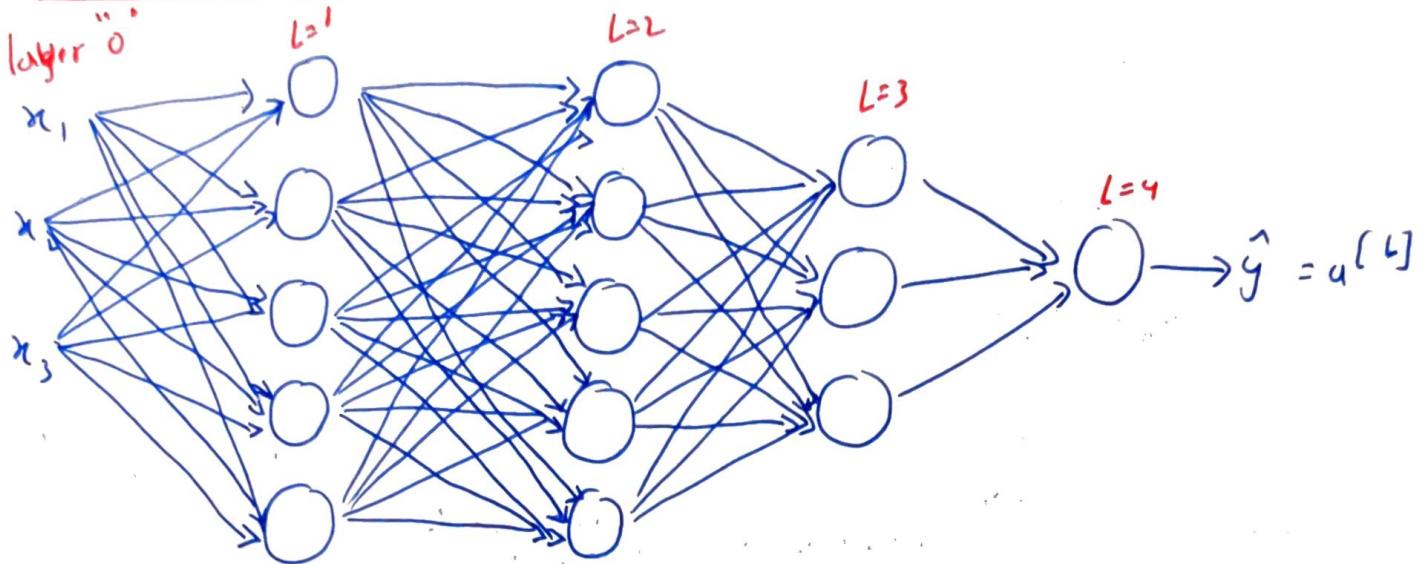
$$b^{[1]} = \text{np.zeros}(2,1)$$

$$\omega^{[2]} = \text{np.random.rand}(1,2) * 0.01$$

$$b^{[2]} = 0$$

$\left. \begin{array}{l} \because \text{should start with very small values} \\ \text{if it is very big, the values after} \\ \text{activat' will be either very large} \\ \text{or very small} \\ \text{making learning slowdown} \\ (\because \text{if sigmoid, on the flatter side of} \\ \text{graph}) \end{array} \right\}$

## DEEP NEURAL NETWORK NOTATION:



$$X = a^{[0]} \quad L = 4 \quad (\text{no. of layers})$$

$h^{[l]}$  = no. of units in layer  $l$

$a^{[l]}$  = activation in layer  $l$

$$a^{[l]} = g^{[l]}(z^{[l]}) \quad \omega^{[l]} = \text{weights for } z^{[l]} \\ b^{[l]} = \text{bias for } z^{[l]}$$

## Forward Propagation in Deep Neural Network:

$$X: z^{[1]} = \omega^{[1]} X + b^{[1]} \quad \xrightarrow{a^{[1]}}$$

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

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

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

$$\vdots \quad \vdots \quad \ddots \quad \vdots$$

$$z^{[4]} = \omega^{[4]} a^{[3]} + b^{[4]}$$

$$a^{[4]} = g^{[4]}(z^{[4]}) = \hat{y}$$

Generalised:  $z^{[l]} = \omega^{[l]} a^{[l-1]} + b^{[l]}$   
 $a^{[l]} = g^{[l]}(z^{[l]})$

Vectorized:  $\begin{aligned} z^{[1]} &= \omega^{[1]} X + b^{[1]} \quad \xrightarrow{A^{[0]}} \\ A^{[1]} &= g^{[1]}(z^{[1]}) \\ z^{[2]} &= \omega^{[2]} A^{[1]} + b^{[2]} \\ A^{[2]} &= g^{[2]}(z^{[2]}) \\ &\vdots \\ \hat{y} &= g^{[4]}(z^{[4]}) = A^{[4]} \end{aligned}$

↑ for loop

## Getting your Matrix dimensions right:

Parameters  $w^{[1]}$  and  $b^{[1]}$        $L = 5$ ,  $n^{[1]} = 3$ ,  $n^{[2]} = 5$ ,  $n^{[3]} = 4$ ,  $n^{[4]} = 2$ ,  $n^{[5]} = 1$

$$z^{[1]} = w^{[1]} \cdot x + b^{[1]}$$

$$\begin{matrix} (3,1) & (3,2) & (4,1) \\ (n^{[1]}, 1) & (n^{[2]}, n^{[1]}) & (n^{[3]}, 1) \end{matrix}$$

$\Rightarrow w^{[1]}$  dimensions:  $(n^{[0]}, n^{[1]})$ , only  $w^{[2]}: (5, 3) : (n^{[2]}, n^{[1]})$

Generally:  $w^{[L]}$  dimensions:  $(n^{[L]}, n^{[L-1]})$

$$\Rightarrow z^{[1]} = \underbrace{w^{[1]} \cdot x}_{(3,1)} + b^{[1]} \Rightarrow (3,1)$$

$$(n^{[1]}, 1) \quad (n^{[0]}, 1)$$

$\Rightarrow$  Generally  $b^{[L]}: (n^{[L]}, 1)$

So when backpropogate:

$$\delta w^{[L]} = w^{[L]} \Rightarrow : (n^{[L]}, n^{[L-1]})$$

$$\delta b^{[L]} = b^{[L]} \Rightarrow : (n^{[L]}, 1)$$

## Vectorized Implementation:

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

$$\Rightarrow z^{[1]} = w^{[1]} \cdot x + b^{[1]}$$

$$(n^{[1]}, m) \quad (n^{[1]}, n^{[0]}) \quad (n^{[0]}, m) \quad (n^{[0]}, 1)$$

: generally:  $z^{[L]}, A^{[L]}: (n^{[L]}, m)$

$$A^{[0]} = x \begin{bmatrix} A_{n^{[0]}}, m \end{bmatrix}$$

$$\delta z^{[L]}, \delta A^{[L]}: (n^{[L]}, m)$$

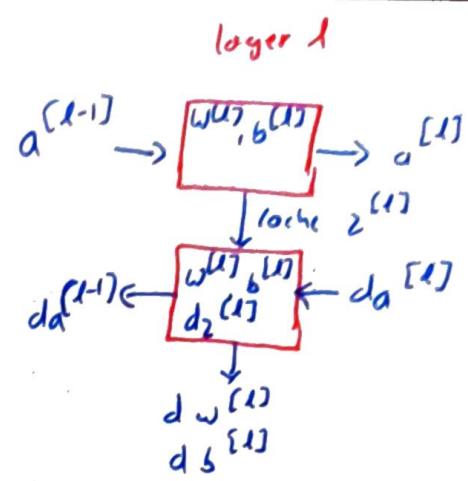
Forward & backward functions

layer  $l$ :  $w^{[l]}, b^{[l]}$

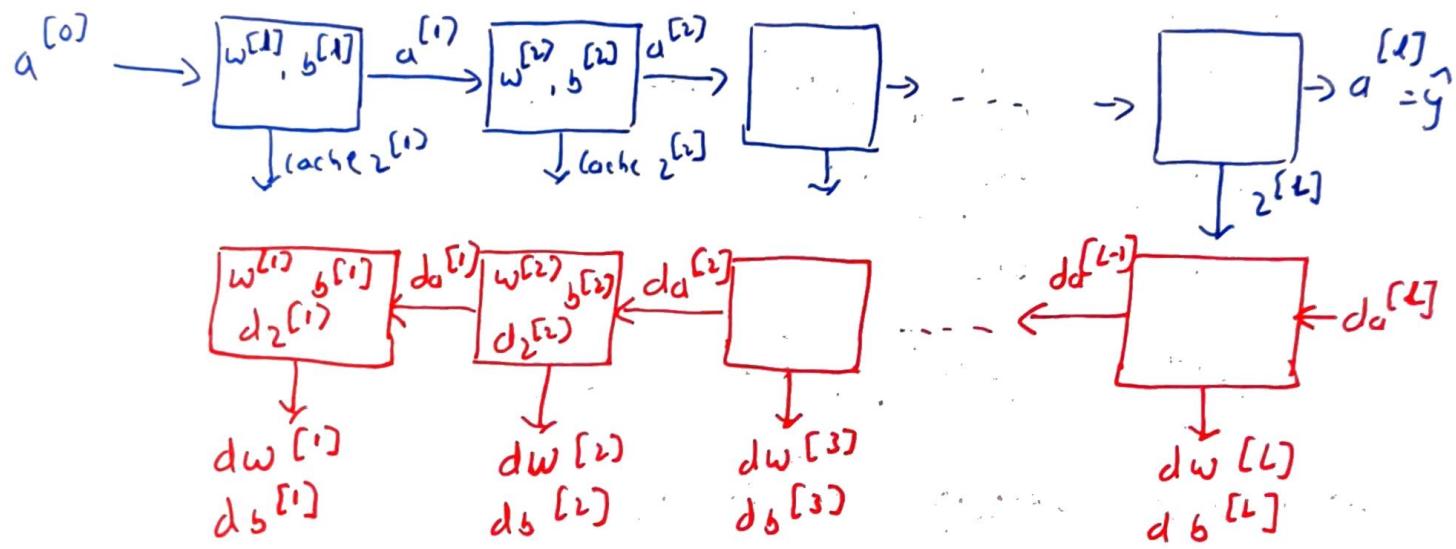
Forward: Input  $a^{[l-1]}$ , Output  $a^{[l]}$

$$z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]} \quad \text{Cache: } z^{[l]}$$

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



Backward: Input:  $da^{[l]}$  Cache  $z^{[l]}$  Output:  $da^{[l-1]}$ ,  $d_w^{[l]}$ ,  $d_b^{[l]}$



$\Rightarrow$  Update:  $w^{[l]} = w^{[l]} - \alpha d_w^{[l]}$

$$b^{[l]} = b^{[l]} - \alpha d_b^{[l]}$$

Forward propagat for layer  $l$ :

Input:  $a^{[l-1]}$   $\rightarrow w^{[l]}, b^{[l]}$

Output:  $a^{[l]}$ , cache ( $z^{[l]}$ )

$$z^{[l]} = w^{[l]} \cdot A^{[l-1]} + b^{[l]}$$

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

Backward propagat:

Input :  $da^{[l]}$

Output :  $da^{[l-1]}, dw^{[l]}, db^{[l]}$

$$dz^{[l]} = \underline{da}^{[l]} * g^{[l]}'(z^{[l]}) \quad \textcircled{1}$$

$$dw^{[l]} = dz^{[l]} \cdot a^{[l-1]}$$

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

$$\underline{da}^{[l-1]} = w^{[l]T} \cdot dz^{[l]} \quad \textcircled{2}$$

Using \textcircled{1} & \textcircled{2} [plugging  $da$  in \textcircled{1}]

$$\Rightarrow dz^{[l]} = w^{[l+1]T} \cdot dz^{[l+1]} * g^{[l]}'(z^{[l]})$$

vectorized (backward):

$$dz^{[l]} = dA^{[l]} * g^{[l]}'(z^{[l]})$$

$$dw^{[l]} = \frac{1}{m} dz^{[l]} \cdot A^{[l-1]T}$$

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

$$dA^{[l-1]} = w^{[l]T} \cdot dz^{[l]}$$

$$dA^{[L]} = \left( -\frac{y^{(1)}}{a^{(1)}} + \frac{(1-y^{(1)})}{(1-a^{(1)})} \quad \dots \quad -\frac{y^{(m)}}{a^{(m)}} + \frac{(1-y^{(m)})}{(1-a^{(m)})} \right) \quad \text{for last step}$$

## PARAMETERS VS HYPERPARAMETERS

Parameters:  $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, w^{[3]}, b^{[3]}, \dots$

Hyperparameters:

- Learning rate  $\rightarrow \alpha$
- Iterations
- # no. of hidden layers  $\rightarrow L$
- no. of hidden units  $\rightarrow n^{[1]}, n^{[2]}, \dots$
- Choice of activation function
- Momentum, mini-batch size, regularization