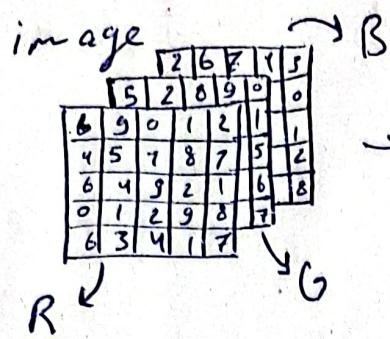


## Binary Classification (Log reg)



→ an image consists of 3 layers (RGB)

Devectorizing the image is to take the matrix and turn it into a vector

$$X = \left[ \underbrace{69_0 \ 12457 \ \dots \ 5289_0}_{R} \ \dots \ \underbrace{2674}_{B} \right] \cdot T$$

The vector size is  $A \times B \times C = 64 \times 64 \times 64$

$$\begin{aligned} \text{the vector size} &= \text{Length of the image} \times \text{width} \times 3 \\ &= 64 \times 64 \times 3 = 12288 \end{aligned}$$

$$nx = 12288$$

## Notation

First training example:-

$$X^{(1)} = [216 \quad 228 \quad 172 \quad 169 \quad \dots \quad \dots \quad \dots]$$

$$y^{(1)} = 0, 1$$

Feature Matrix  $X$ : Contains all feature vectors  $x$

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

The height is the number of features.

The width is the number of training examples.

Lable vector  $\mathbf{Y}_x$

$$\mathbf{Y} = [y^{(1)} \ y^{(2)} \ y^{(3)} \ \dots \ \dots \ y^{(m)}]$$

shape of  $X \rightarrow (n \times m)$

shape of  $\mathbf{Y} \rightarrow (1, m)$

Log reg cost function:

For log reg we don't use a simple loss function

Like  $L(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2$  because the curve of the  $w, b$  to the loss will be non convex

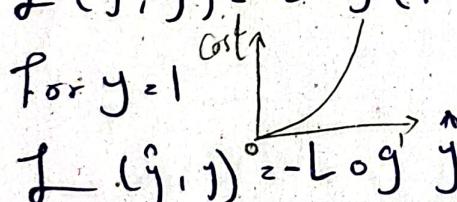


instead we us something more convex  
like the log loss function

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

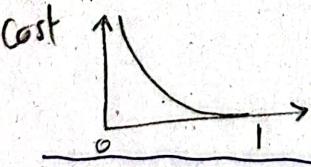
For  $y=0$

$L(\hat{y}, y) = -\log(1-\hat{y}) \rightarrow$  for the loss to be 0 you need  $-\log(1-\hat{y})$  to be zero so you need  $\hat{y}$  to be as small as possible  $\approx 0$



For  $y=1$

$L(\hat{y}, y) = -\log \hat{y} \rightarrow$  for the loss to be 0 you need  $-\log \hat{y}$  to be 0 and since  $-\log(1)=0$  you need  $\hat{y}$  to be as big as possible  $\approx 1$



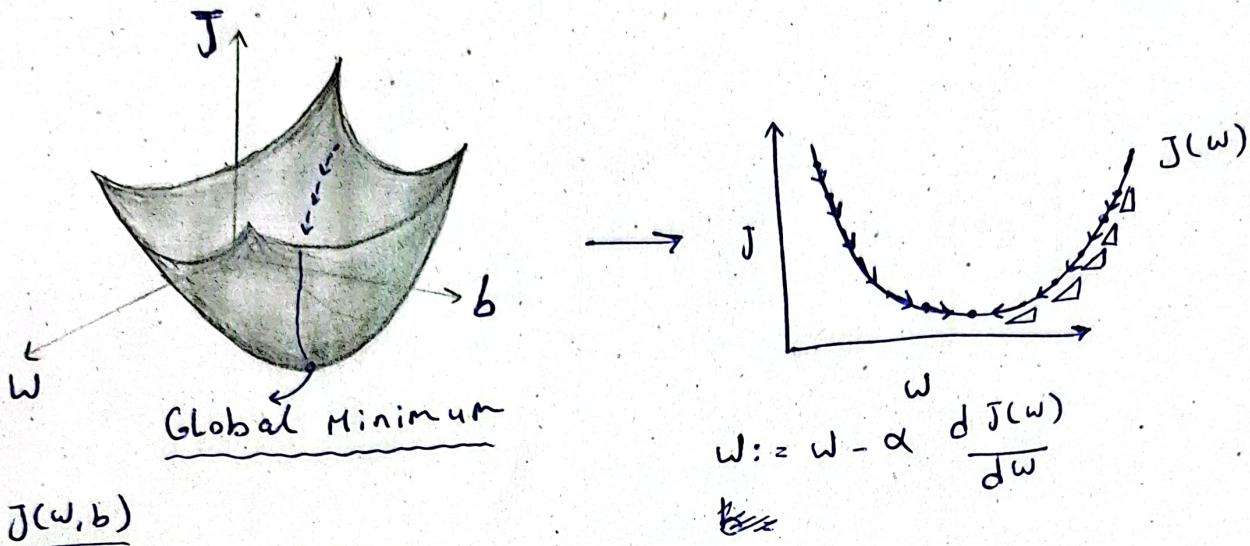
\* Loss function works for a single training example  
but to measure the loss for the whole dataset we'll define a Cost function.

Cost function:  $J(w, b) = \frac{1}{m} \sum_{i=1}^m -[y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log(1-\hat{y}^{(i)})]$

num of training examples

## Gradient descent

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}_i, y_i) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)})$$

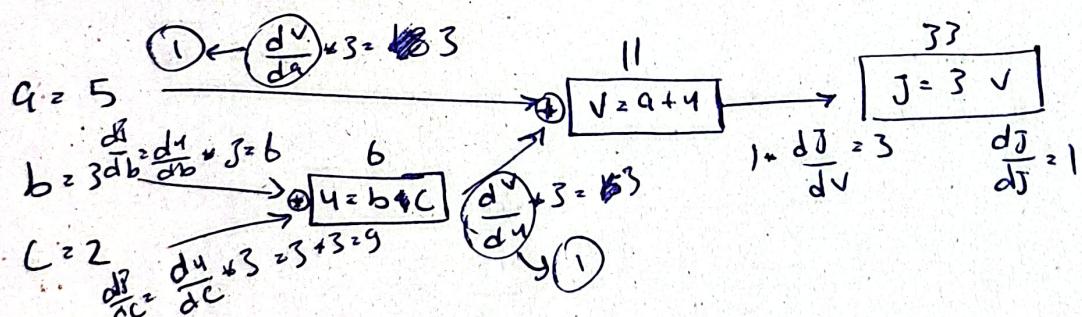


$J(w, b)$

$$\begin{aligned} w &:= w - \alpha \frac{dJ}{dw} \\ b &:= b - \alpha \frac{dJ}{db} \end{aligned} \quad \text{in coding}$$

## Computation graph

$$J(a, b, c) = 3(a + bc) \quad \text{if } u = a + bc, \quad V = a + u$$



$$\frac{dJ}{dv} = 3(a + u) = 3$$

$$\text{if } V = a + u \rightarrow \frac{dv}{du} = a, \quad \frac{dv}{da} = u$$

$$\left| \begin{array}{c} \frac{dJ}{da} \quad \frac{dJ}{db} \quad \frac{dJ}{dc} \end{array} \right|$$

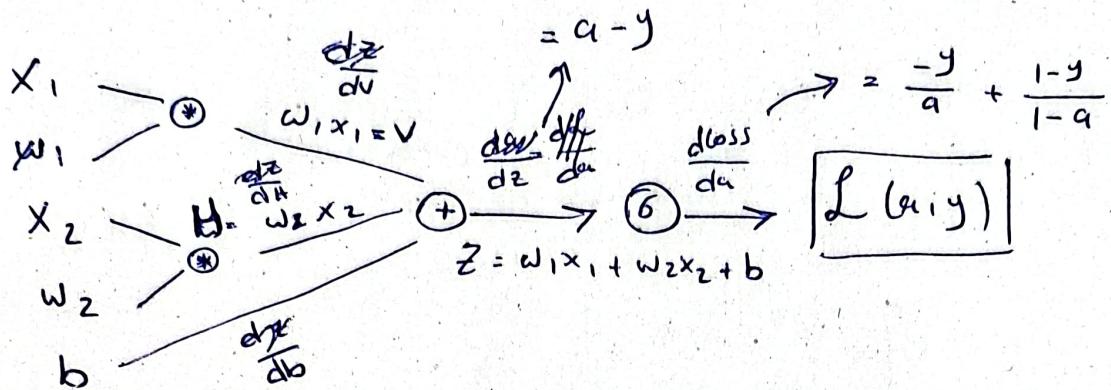
$$\begin{aligned} \frac{dJ}{da} &= 1, \quad \frac{dJ}{dv} = \frac{dJ}{dJ} \frac{dJ}{dv} = 1 * 3 = 3 \\ \frac{dJ}{du} &= \frac{dJ}{dv} \frac{dv}{du} = 3 * 4 = 12 \\ \frac{dJ}{da} &= \frac{dJ}{dv} \frac{dJ}{du} = 3 * 4 = 12 \end{aligned}$$

$$\begin{aligned} \frac{dJ}{du} &= 3 \\ \frac{dJ}{du} &= 9 \\ \frac{dJ}{db} &= \frac{dJ}{dv} \frac{dv}{du} \frac{du}{db} = 3 * 1 * 2 = 6 \end{aligned}$$

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

$$\hat{y} = a = \sigma(Z)$$

$$L(a, y) = - (y \log(a) + (1-y) \log(1-a))$$



$$da = \frac{dL}{da}$$

$$dz = \frac{dL}{da} \frac{da}{dz} = \boxed{da \cdot \frac{da}{dz}}$$

$$dv = \frac{dL}{da} \frac{da}{dz} \frac{dz}{dv} = \boxed{dz \cdot \frac{dz}{dv}}$$

$$dH = \frac{dL}{da} \frac{da}{dz} \frac{dz}{dH} = \boxed{dz \cdot \frac{dz}{dH}}$$

$$db = \frac{dL}{da} \frac{da}{dz} \frac{dz}{db} = \boxed{dz \cdot \frac{dz}{db}} = dz$$

$$dx_1 = \frac{dL}{da} \frac{da}{dz} \frac{dz}{dv} \frac{dv}{dx_1} = \boxed{d.v \cdot \frac{dv}{dx_1}} = \underline{\underline{dz + w_1}}$$

$$dw_1 = \frac{dL}{da} \frac{da}{dz} \frac{dz}{dv} \frac{dv}{dw_1} = \boxed{d.v \frac{dv}{dw_1}} = \underline{\underline{dz * 1 * x_1 - dz \cdot x_1}}$$

$$dx_2 = \frac{dL}{da} \frac{da}{dz} \frac{dz}{dv} \frac{dv}{dx_2} = \boxed{d.H \cdot \frac{dH}{dx_2}} = \underline{\underline{dz \cdot w_2}}$$

$$dw_2 = \frac{dL}{da} \frac{da}{dz} \frac{dz}{dv} \frac{dv}{dw_2} = \boxed{d.H \cdot \frac{dH}{w_2}} = \underline{\underline{dz \cdot x_2}}$$

$$\cancel{db = \frac{dL}{da} \frac{da}{dz} \frac{dz}{dv} \frac{dv}{db}}$$

$$\bar{J}(\omega, b) = \frac{1}{m} \sum_{i=1}^m L(a^{(i)}, y^{(i)})$$

$$\rightarrow a^{(i)} = g^{(i)} = \sigma(z^{(i)}) = \sigma(w^\top x^{(i)} + b) \Leftrightarrow [\underline{d\omega_1^{(i)}}, \underline{d\omega_2^{(i)}}, \underline{db^{(i)}}]$$

$$\frac{d}{d\omega_1} J(\omega, b) = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial \omega_1} L(a^{(i)}, y^{(i)})$$

$$= \frac{d\ell}{da} \frac{da}{dz} \frac{dz}{dv} \frac{dv}{d\omega_1} = \boxed{d\omega_1}$$

$$J = 0, d\omega_1 = 0, d\omega_2 = 0, db = 0$$

For  $i = 1$  to  $m$

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

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

$$\bar{J} += \left[ j^{(i)} \log(a^{(i)}) + (1-y^{(i)}) \log(1-a^{(i)}) \right]$$

$$dz^{(i)} = a^{(i)} - y^{(i)}$$

$$d\omega_1 += x_1^{(i)} dz^{(i)}$$

$$d\omega_2 += x_2^{(i)} dz^{(i)}$$

$$db += dz^{(i)}$$

$$J / = m$$

$$d\omega_1 / = m$$

$$d\omega_2 / = m$$

$$db / = m$$

after the ~~first~~ loop

$$d\omega_1 = \frac{dJ}{d\omega_1}$$

so  $d\omega_1$  is the accumulation of derivatives across the training set.

$$w_1 := w_1 - \alpha d\omega_1$$

$$w_2 := w_2 - \alpha d\omega_2$$

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

each step

## Vectorizing logistic regression

$$Z^{(1)} = \mathbf{w}^T \cdot \mathbf{x}^{(1)} + b, Z^{(2)} = \mathbf{w}^T \cdot \mathbf{x}^{(2)} + b, Z^{(3)} = \mathbf{w}^T \cdot \mathbf{x}^{(3)} + b$$

$$a^{(1)} = \sigma(Z^{(1)}), a^{(2)} = \sigma(Z^{(2)}), a^{(3)} = \sigma(Z^{(3)})$$

$$\mathbf{X} = \left[ \begin{array}{c|c|c|c|c} & & & & \\ \mathbf{x}^{(1)} & \mathbf{x}^{(2)} & \dots & \mathbf{x}^{(n)} & \\ & & & & \end{array} \right] \quad \begin{matrix} \uparrow \\ n \rightarrow \text{number of training examples.} \end{matrix}$$

↳ number of features.

$$\mathbf{X}.\text{shape} = (n \times m)$$

$$\begin{aligned} Z &= [Z^{(1)} \ Z^{(2)} \ Z^{(3)} \ Z^{(m)}] = \mathbf{w}^T \cdot \mathbf{X} + [b, b_2, b_3, \dots, b_m]^T \\ &\Rightarrow \mathbf{w}^T \cdot \left[ \begin{array}{c|c|c|c|c} & & & & \\ \mathbf{x}^{(1)} & \mathbf{x}^{(2)} & \dots & \mathbf{x}^{(n)} & \\ & & & & \end{array} \right]_{n \times m} \\ &= \boxed{\mathbf{w}^T \cdot \mathbf{x}^{(1)} + b} \quad \boxed{\mathbf{w}^T \cdot \mathbf{x}^{(2)} + b} \quad \dots \quad \boxed{\mathbf{w}^T \cdot \mathbf{x}^{(n)} + b} \\ &\quad Z^{(1)} \qquad \qquad \qquad Z^{(2)} \qquad \qquad \qquad Z^{(m)} \end{aligned}$$

$$Z = np.\text{dot}(\mathbf{w}^T \cdot \mathbf{X}) + b \quad \underbrace{\text{python broadcasting.}}$$

$$A = \sigma(Z) = [\sigma(Z^{(1)}) \ \sigma(Z^{(2)}) \ \dots \ \sigma(Z^{(m)})]$$

$$dZ = A - Y \quad (1 \times m)$$

$$db = \frac{1}{n} np.\text{sum}(dZ)$$

$$dw = \frac{1}{n} X \cdot dZ^T$$

$$Z = np.\text{dot}(\mathbf{w}^T, \mathbf{X}) + b$$

$$A = \sigma(Z)$$

$$dZ = A - Y$$

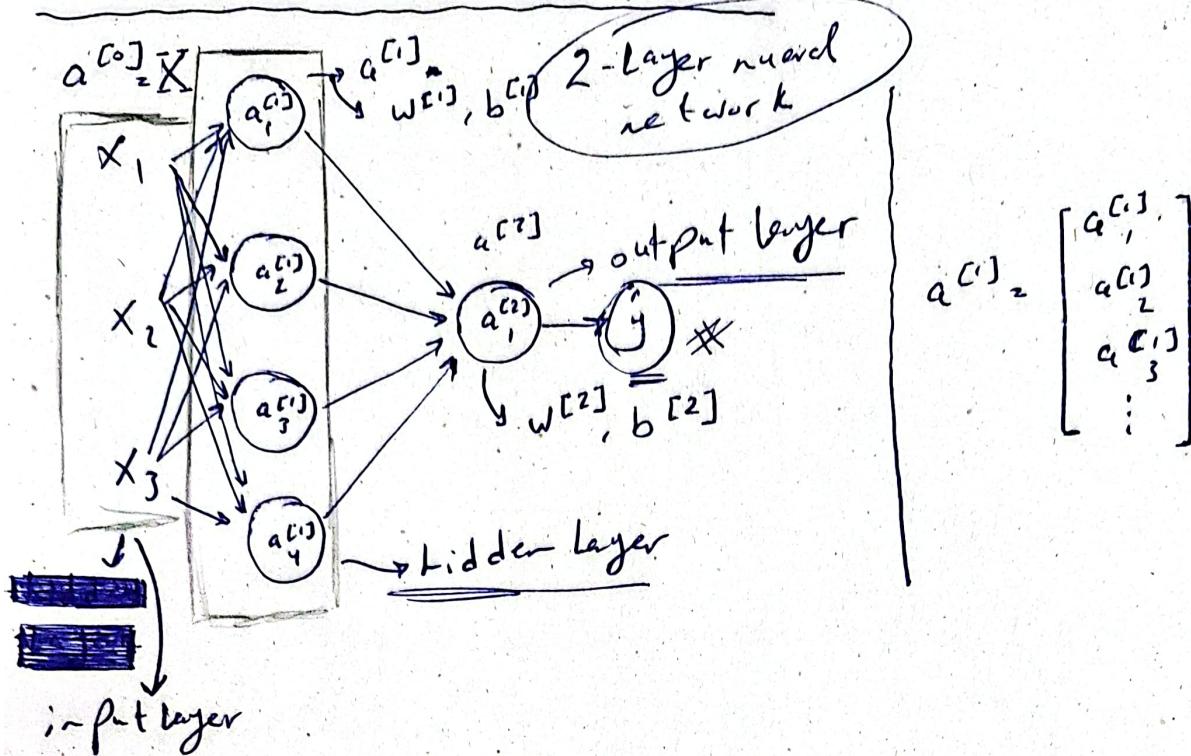
$$dw = \frac{1}{n} X \cdot dZ^T$$

$$db = \frac{1}{n} np.\text{sum}(dZ)$$

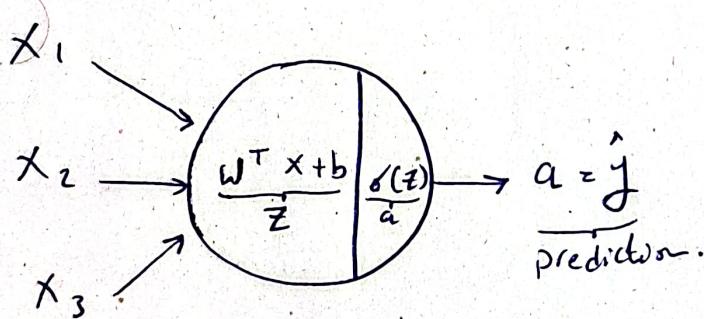
$$w_i = w_i - \alpha dw \rightarrow *$$

$$b_i = b_i - \alpha db \rightarrow *$$

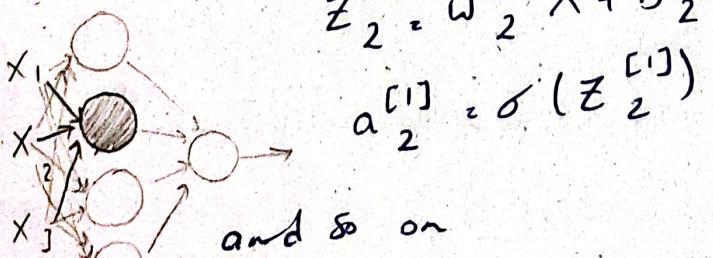
# Neural network representation



~~$a^{[0]}$~~  Computing the output of a neural network

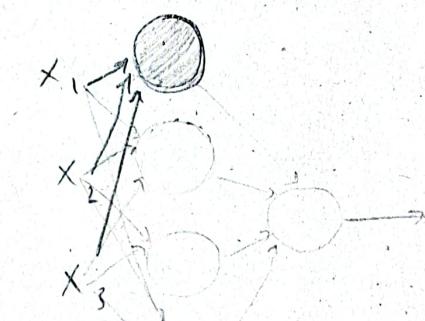


$$z = w^T \cdot x + b \rightarrow a = \sigma(z)$$



$$Z_i^{(1)} = w_i^{(1)T} x + b_i^{(1)}$$

$$a_i^{(1)} = \sigma(Z_i^{(1)})$$



$$Z_1^{(1)} = w_1^{(1)T} x + b_1^{(1)}$$

$$a_1^{(1)} = \sigma(Z_1^{(1)})$$

general for -

$$Z_{(i)}^{(j)} = w_{(i)}^{(j)T} x + b_{(i)}^{(j)}$$

$$a_{(i)}^{(j)} = \sigma(Z_{(i)}^{(j)})$$

→ in the previous Logistic regression as a single neuron example:

- \*  $X$  was a matrix  $X$  consisting of all images stacked horizontally.

$$\bar{X} = \begin{bmatrix} 1 & | & | & | \\ x_{(1)} & x_{(2)} & \dots & x_{(m)} \\ 1 & | & | & | \end{bmatrix} \quad \begin{matrix} \uparrow \\ n \times p \end{matrix}$$

$\leftarrow m$

and  $w$  was a row vector after it was transposed  
so each feature has a weight

$$Z = w^T \cdot \bar{X} + b$$

$$= [w_1 \ w_2 \ w_3 \ w_4 \dots] \begin{bmatrix} 1 & | & | & | \\ x_{(1)} & x_{(2)} & \dots & - \\ 1 & | & | & | \end{bmatrix} + b$$

↳ scalar

\* each weight is going to be multiplied with a specific feature resulting in their dot product so we can calculate the overall cost.

$$Z = w^T \cdot \bar{X} + b \rightarrow A = \sigma(Z) \rightarrow \boxed{j = \sum_{i=1}^n L(y_i, A)}$$

\* When handling a whole layer then  $w$  is going to be multiple vectors - each corresponding to a single neuron - stacked together in a matrix, and  $b$  is going

$$W^T = \begin{bmatrix} \underline{w_1^{[1]}} \\ \underline{w_2^{[1]}} \\ \underline{w_3^{[1]}} \\ \vdots \\ \underline{w_n^{[1]}} \end{bmatrix} \quad \text{to be a vector.}$$

$$b = [b_1^{[1]}, b_2^{[1]}, \dots, b_n^{[1]}]$$

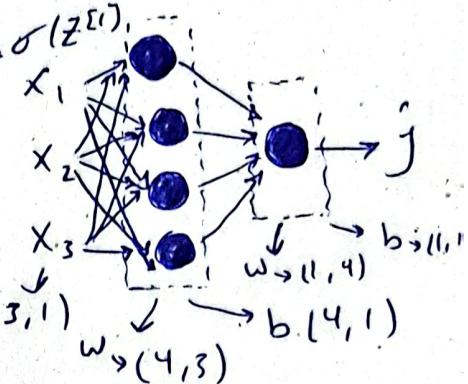
$$Z^{[1]} = \begin{bmatrix} w_1^{[1]T} \\ w_2^{[1]T} \\ \vdots \\ w_n^{[1]T} \end{bmatrix} \cdot \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_{(1)} & x_{(2)} & \dots & x_{(m)} \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ \vdots \\ b_n^{[1]} \end{bmatrix}$$

$$Z^{[1]} = w^{[1]T} \cdot X + b \Rightarrow A^{[1]} = \sigma(Z^{[1]})$$

(4,1)      (4,3)      (3,1)      (4,1)

$$Z^{[2]} = w^{[2]T} \cdot A^{[1]} + b$$

(1,1)      (1,4)      (4,1)      (1,1)



$$\downarrow A^{[2]} = \sigma(Z^{[2]}) = \hat{y}$$

For multiple training examples

$$X = \begin{bmatrix} 1 & 1 & 1 & \dots & 1 \\ x_{(1)} & x_{(2)} & \dots & x_{(m)} & \end{bmatrix}_{n \times p}^T$$

$$\begin{aligned} Z^{[1]} &= w^{[1]T} X + b^{[1]} \\ A^{[1]} &= \sigma(Z^{[1]}) \\ Z^{[2]} &= w^{[2]T} \cdot A^{[1]} + b^{[2]} \\ A^{[2]} &= \hat{y} = \sigma(Z^{[2]}) \end{aligned}$$

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

$$4 \times 4 \times 3 = 16 \times 3 = 48$$

**Important**

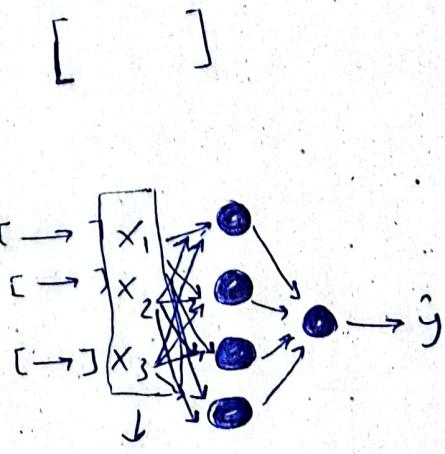
$$Z^{[1]} = w^{[1]T} \cdot X + b^{[1]} \rightarrow A^{[1]} = \sigma(Z^{[1]})$$

(4,10)      (4,48)      (48,10)      (4,10)

$$Z^{[2]} = w^{[2]T} \cdot A^{[1]} + b^{[2]}$$

(1,10)      (1,4)      (4,10)      (1,10)

↳ output for 10-images prediction.



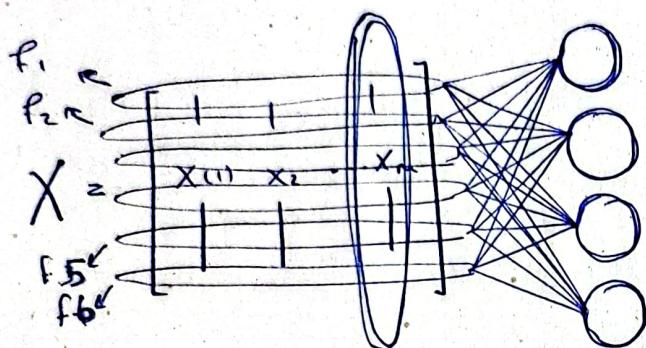
$$x = \begin{bmatrix} x_1 & x_2 & \dots & x_m \end{bmatrix}^T$$

Justification for the vectorized implementation

① .

$$\begin{aligned} Z^{(1)} &= w^{(1)T} \cdot x^{(1)} + b^{(1)} \\ Z^{(2)} &= w^{(2)T} \cdot x^{(2)} + b^{(2)} \\ Z^{(3)} &= w^{(3)T} \cdot x^{(3)} + b^{(3)} \end{aligned}$$

$$w^{[1]} = \left[ \begin{array}{c} \hline \\ \hline \\ \hline \\ \hline \end{array} \right]$$



→ Let's imagine that we do this process - one image at a time -.

\* First image has 6 features and each feature is going to be an input to 4 neurons of the first layer.

$$x_1 \rightarrow (b, 1), \quad w^{[1]} \rightarrow (4, b)$$

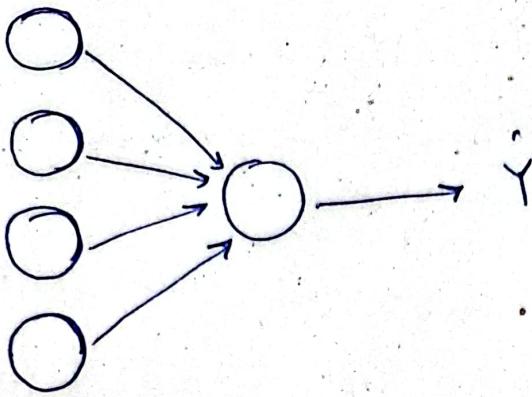
$$\cancel{Z^{(1)}} \quad w^{(1)}_{x^{(1)}} = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix} \quad w^{(1)}_{x^{(2)}} = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}$$

$$Z^{[E]}_z = \left[ \begin{array}{c|c|c|c} \vdots & \vdots & \vdots & \vdots \\ \hline & & & \end{array} \right] \quad \begin{matrix} n \\ \downarrow y \end{matrix}$$

$$\omega^{[1]T} x^{(1)} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \xrightarrow{\text{neuron 1}} \dots \xrightarrow{\text{neuron 2}} \dots \xrightarrow{\text{--- 3}} \dots \xrightarrow{\text{--- 4}}$$

→ these are the activations of the 4 neurons on the 6 features of the first training example

(2)



$$X = \begin{bmatrix} & & & & 1 \\ x_{(1)} & \dots & \dots & x_{(10)} & \\ & & & & 1 \\ 1 & & & & 1 \\ \downarrow & & & & \downarrow \\ (6, 10) & & & & (4, 10) \end{bmatrix} \Rightarrow Z^{[1]} = \begin{bmatrix} 1 \\ \vdots \\ \vdots \\ \vdots \\ \vdots \\ 1 \end{bmatrix}$$

$$\omega^{[1]} = (4, 6)$$

$$\boxed{A^{[1]} = \sigma(Z^{[1]}) \rightarrow (4, 10)}$$

$$Z^{[2]} = \omega^{[2]T} \cdot A^{[1]} + b^{[2]}$$

$$Z^{[2](1)} = \omega^{[2]} \cdot A^{[1](1)}$$

$$Z^{[2](2)} = \omega^{[2]} \cdot A^{[1](2)}$$

$$Z^{[2](1)} = [ \cdot ] , Z^{[2](2)} = [ \cdot ] , Z^{[2](3)} = [ \cdot ]$$

$$Z^{[2]} = \begin{bmatrix} \cdot & \cdot & \cdot & \cdots & \cdot \end{bmatrix}$$

$$A^{[2]}, Y = \underline{\underline{b(Z^{[2]})}}$$

↳ predictions for the classes of  
the 10 images

## Gradient descent for 1-hidden layer network

Params  $\rightarrow w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}$

dims  $\rightarrow n_x = \frac{n^{[0]}}{T_x}, n^{[1]}, n^{[2]} = 1$   
 $n^{[1]}$   $\hookrightarrow$  hidden units

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

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

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

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

Cost  $= J(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]})$

$$= \frac{1}{m} \sum_{i=1}^m \ell(y_i, A^{[2]} = \hat{y}_i)$$

GD

↳ Repeat

$\rightarrow$  Compute the Predictions ( $\hat{y}^{(1)}, \dots, \hat{y}^{(m)}$ )

$\rightarrow$  ~~Compute the Cost~~  $\rightarrow$  Forward prob

$\rightarrow$  Compute the derivatives  $d w^{[1]} = \frac{d J}{d w^{[1]}}, d b^{[1]} = \frac{d J}{d b^{[1]}}$

$d w^{[2]}, d b^{[2]}$

$\Rightarrow$  back prob

$\rightarrow$  update  $\rightarrow w^{[1]} = w^{[1]} - \alpha d w^{[1]}$

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

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

$b^{[2]} = b^{[2]} - \alpha d b^{[2]}$

Forward prob

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

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

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

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

back prob

$$d z^{[2]} = A^{[2]} - Y$$

$$d w^{[2]} = \frac{1}{m} d z^{[2]} A^{[1]T}$$

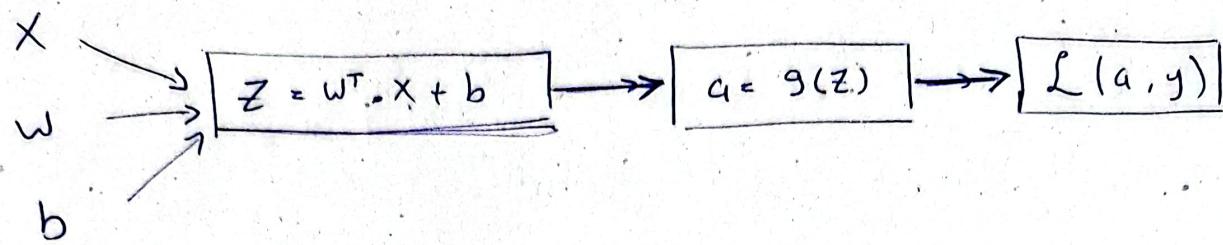
$$d b^{[2]} = \frac{1}{m} \text{np.sum}(d z^{[2]})$$

$$d z^{[1]} = \frac{1}{m} d w^{[2]} T d z^{[2]} * g^{[2]}'(z^{[2]})$$

$$d w^{[1]} = \frac{1}{m} d z^{[1]} A^{[0]T}$$

$$d b^{[1]} = \frac{1}{m} \text{np.sum}(d z^{[1]})$$

## Back propagation intuition



Forward prob

$$Z = w^T \cdot X + b \rightarrow a = g(Z) \rightarrow L(a, y) = -y \log a - (1-y) \log(1-a)$$

Cost ↓

Back prob

$$\frac{dL}{dL} = 1$$

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

"da"

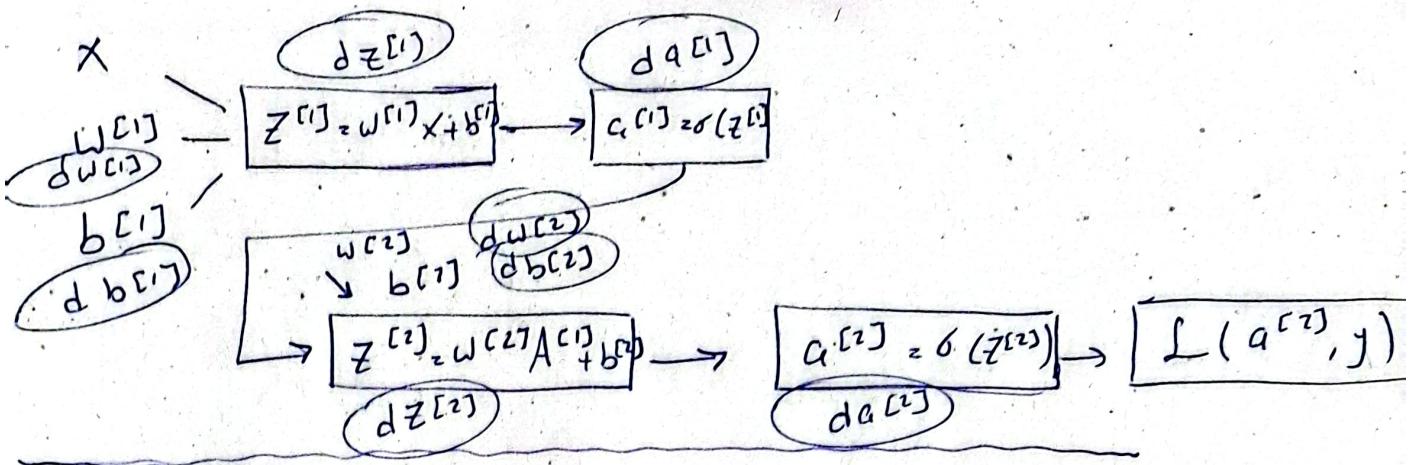
$$dZ = a - y = da \cdot g'(Z) \Rightarrow g(Z) = \sigma(Z)$$

$$dZ = \frac{dL}{dZ} = \left[ \frac{dL}{da} \right] \left[ \frac{da}{dZ} \right] = "da" \cdot \left[ \frac{d}{dZ} (\sigma(Z)) \right] \sigma'(Z)$$

"da" = da · g'(Z)

$$dW = dZ \cdot X$$

$$db = dZ$$



### Forward Prob

$$Z^{[1]} = W^{[1]T} X + b^{[1]} \rightarrow a^{[1]} = \sigma(Z^{[1]})$$

$$Z^{[2]} = W^{[2]T} a^{[1]} + b^{[2]} \rightarrow a^{[2]} = \sigma(Z^{[2]})$$

$$\text{Cost} = -y \log a^{[2]} - (1-y) \log(1-a^{[2]})$$

### Back Prob

$$A = \sigma(Z)$$

$$\frac{da}{dz} = \sigma'(Z) = A(1-A)$$

$$\rightarrow da^{[2]} = \frac{-y}{a^{[2]}} + \frac{1-y}{1-a^{[2]}}$$

$$\rightarrow dz^{[2]} = \frac{dL}{da^{[2]}} \cdot \frac{da^{[2]}}{dz^{[2]}} = \frac{-y}{a^{[2]}} + \frac{1-y}{1-a^{[2]}} * (a^{[2]}(1-a^{[2]}))$$

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

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

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

$$\rightarrow da^{[1]} = \underbrace{\frac{dL}{da^{[2]}}}_{=a^{[2]}-y} \cdot \frac{da^{[2]}}{dz^{[2]}} \cdot \frac{d\cancel{z}^{[2]}}{da^{[1]}} = \cancel{(a^{[2]}-y)} \cancel{\left(\frac{-y}{a^{[2]}} + \frac{1-y}{1-a^{[2]}}\right)}$$

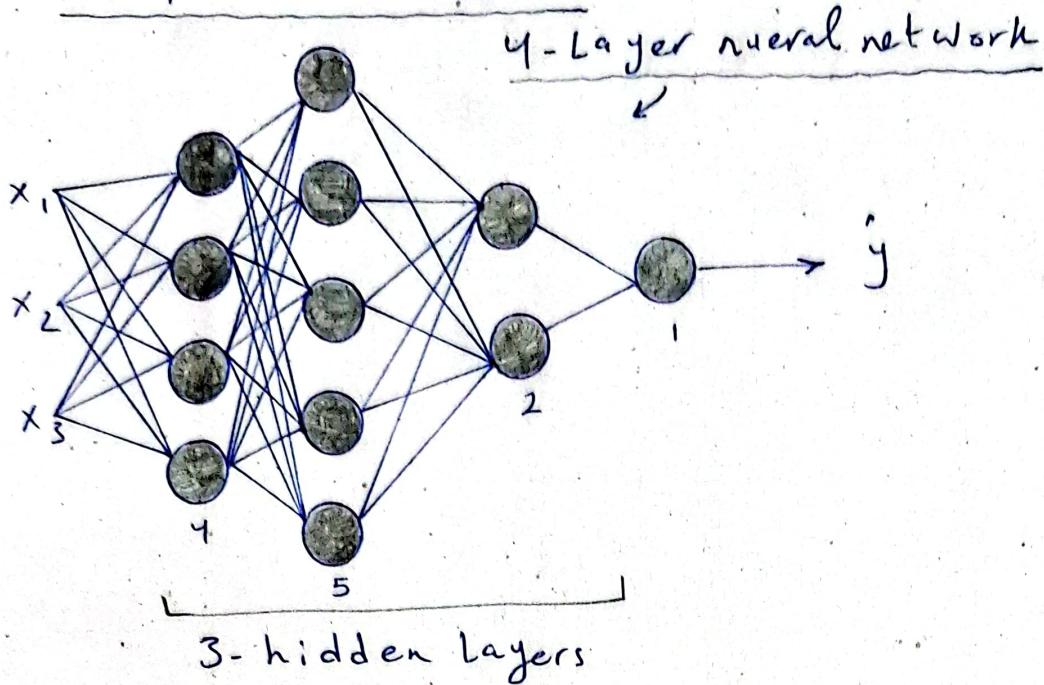
$$= a^{[2]} - y \cdot w^{[2]} = dz^{[2]} \cdot w^{[2]}$$

$$\rightarrow \frac{dz^{[1]}}{da^{[1]}} = \cancel{a^{[1]}} \cdot (1-a^{[1]}) \cancel{g^{[1]}'(z^{[1]})} \cdot z^{[2]}$$

$$\cdot dz^{[1]} = da^{[1]} \cdot \frac{dz^{[1]}}{da^{[1]}} = \cancel{(a^{[2]}-y)} \cancel{w^{[2]}} \cdot$$

$$= dz^{[2]} \cdot w^{[2]} * \cancel{g^{[1]}'(z^{[1]})}$$

# Deep neural networks



$$\rightarrow L = 4 \quad (\text{* Layers})$$

$\rightarrow n^{[l]}$  = units in Layer  $l \rightarrow n^{[1]} = 4, n^{[2]} = 5, 2, 1 \dots$

$\rightarrow n^{[0]} = n_x = 3 \rightarrow \text{input features}$

$\rightarrow a^{[l]}$  = activations in layer  $l$

$\rightarrow a^{[l]} = g^{[l]}(z^{[l]}) \rightarrow a^{[0]} = X, a^{[L]} = \hat{y}$

$\rightarrow w^{[l]}$  = weights for  $z^{[l]}$

$\rightarrow b^{[l]}$  = biases for  $z^{[l]}$

## Forward propagation

- $z^{[1]} = w^{[1]} \cdot X + b^{[1]}$
- $a^{[1]} = g^{[1]}(z^{[1]})$
- $z^{[2]} = w^{[2]} \cdot a^{[1]} + b^{[2]}$
- $a^{[2]} = g^{[2]}(z^{[2]})$
- $z^{[3]} = w^{[3]} \cdot a^{[2]} + b^{[3]}$
- $a^{[3]} = g^{[3]}(z^{[3]})$
- $z^{[4]} = w^{[4]} \cdot a^{[3]} + b^{[4]}$
- $a^{[4]} = g^{[4]}(z^{[4]})$

most likely sigmoid.

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

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

## Vectorized Form

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

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

⋮

⋮

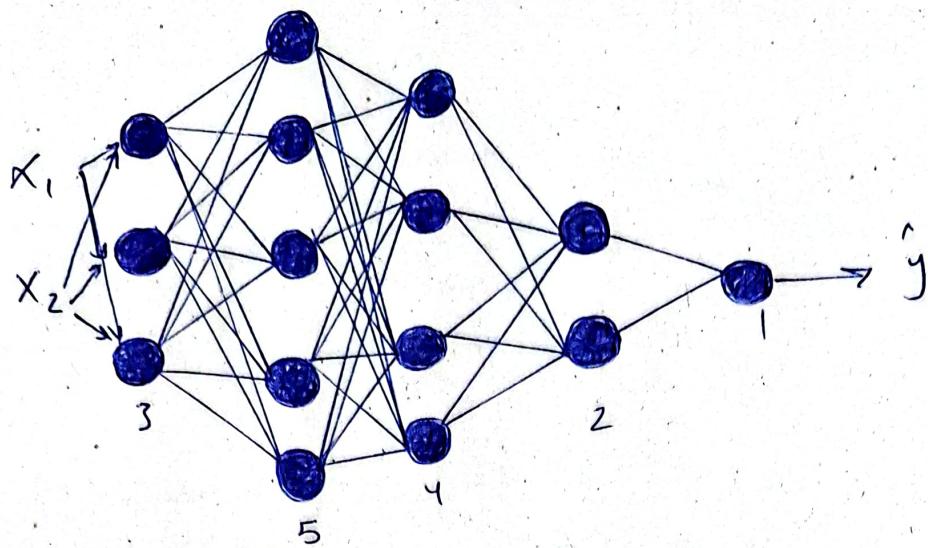
⋮

$$Z^{[4]} = W^{[4]} \cdot A^{[3]} + b^{[4]}$$

$$A^{[4]} = g^{[4]}(Z^{[4]})$$

$\hat{Y} \rightarrow \text{predictions}$

$$X = A^{[0]}$$



$$Z^{[1]} = W^{[1]} \cdot X + b^{[1]}$$

$$(3, 1) \quad (3, 2) \quad (2, 1) \quad (3, 1)$$

$$\downarrow \quad \downarrow \quad \downarrow \quad \downarrow$$

$$(n^{[e]}, m) \quad (n^{[e]}, n^{[e-1]}) \quad (n_x, m)$$

$$W \rightarrow (n^{[e]}, n^{[e-1]})$$

$$b \rightarrow (n^{[e]}, 1)$$

Let's imagine that we have 10 training examples.

$$\rightarrow Z^{[1]} = W^{[1]} \cdot X + b^{[1]}$$

$$(3, 10) \quad (3, 2) \quad (2, 10) \quad (3, 1)$$

$$\rightarrow Z^{[2]} = W^{[2]} \cdot A^{[1]} + b^{[2]}$$

$$(5, 10) \quad (5, 3) \quad (3, 10) \quad (5, 1)$$

$$\rightarrow Z^{[3]} = W^{[3]} \cdot A^{[2]} + b^{[3]}$$

$$(4, 10) \quad (4, 5) \quad (5, 10) \quad (4, 1)$$

$$\rightarrow Z^{[4]} = W^{[4]} \cdot A^{[3]} + b^{[4]}$$

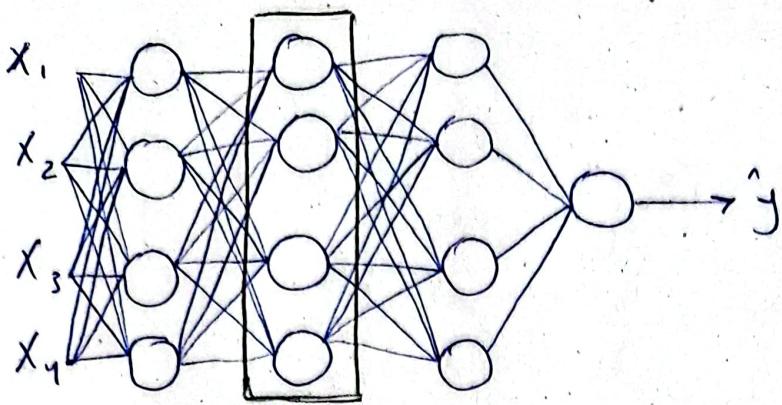
$$(2, 10) \quad (2, 4) \quad (4, 10) \quad (2, 1)$$

$$\rightarrow Z^{[5]} = W^{[5]} \cdot A^{[4]} + b^{[5]}$$

$$(1, 10) \quad (1, 2) \quad (2, 10) \quad (1, 1)$$

$$A^{[5]} = \delta(Z^{[5]}) \rightarrow \text{predictions}$$

b is broadcasted  
by python  
automatically.



Layer 2  $\rightarrow w^{[2]}, b^{[2]}$

$$\rightarrow z^{[2]} = w^{[2]} \cdot A^{[1]} + b^{[2]} \rightarrow A^{[2]} = g^{[2]}(z^{[2]})$$

(input:  $A^{[1]}$ , output:  $z^{[2]}, A^{[2]}$ )

Back prop for Layer 2

~~for forward~~  
input,  $da^{[2]}$ , output:  $da^{[1]}, dw^{[2]}, db^{[2]}$   
+ Cache ( $z^{[2]}$ )

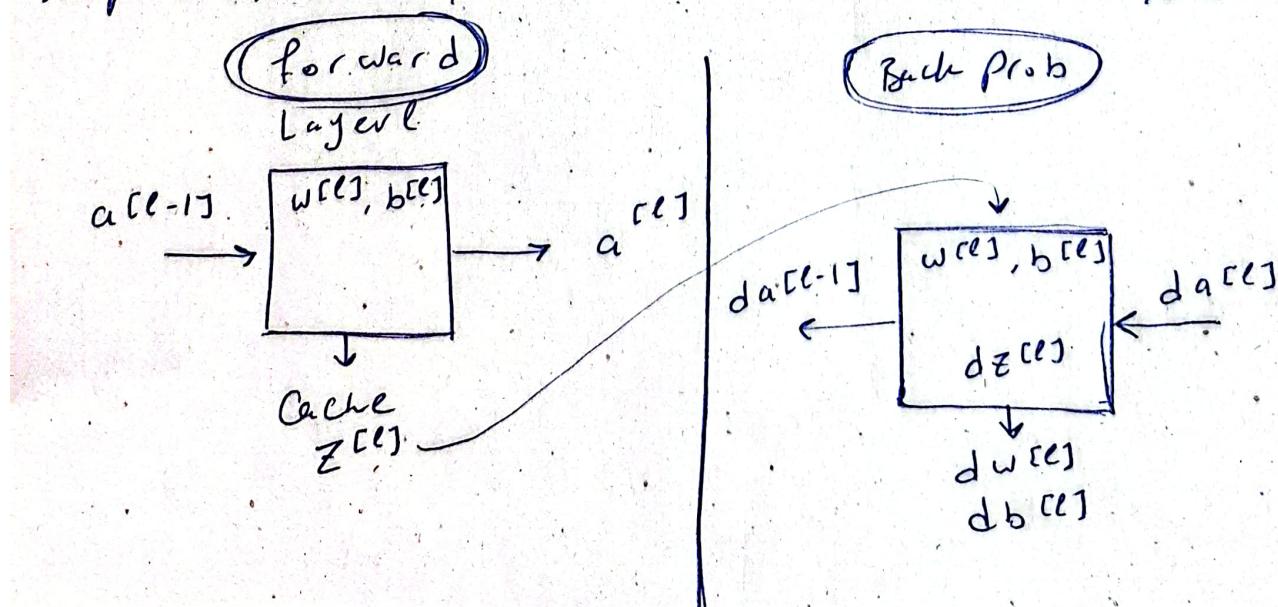
for Layer  $l$ :  $w^{[l]}, b^{[l]}$

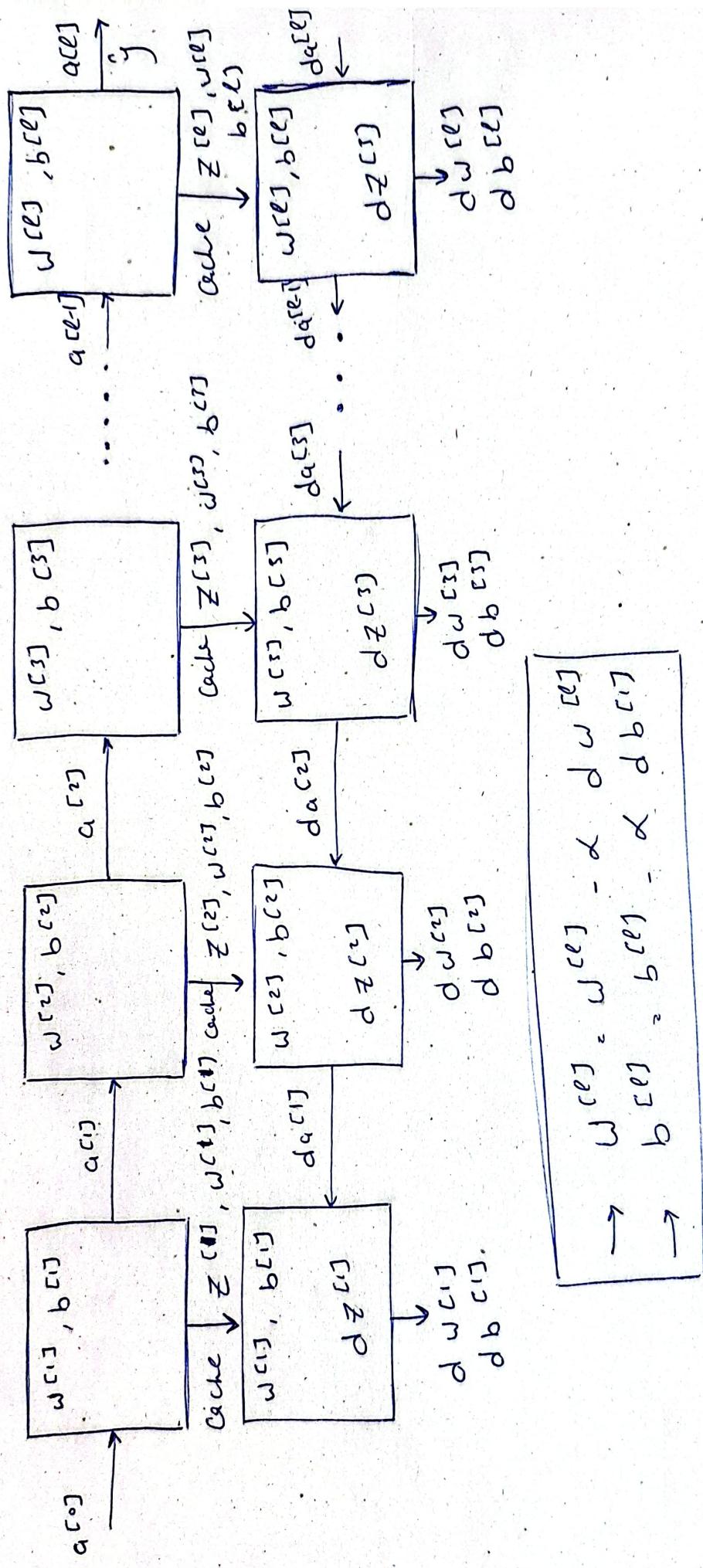
$$\rightarrow z^{[l]} = w^{[l]} \cdot a^{[l-1]} + b^{[l]}$$

for ward  
input  $\rightarrow a^{[l-1]}$ , output  $\rightarrow a^{[l]}$

Back prop

input  $\rightarrow da^{[l]}$ , output  $\rightarrow da^{[l-1]}, dw^{[l]}, db^{[l]}$





## forward propagation

- in put  $a^{[l-1]}$
  - out put  $a^{[l]}$ , cache ( $z^{[l]}$ )
- $$\rightarrow z^{[l]} = w^{[l]} \cdot a^{[l-1]} + b^{[l]}$$
- $$a^{[l]} = g^{[l]}(z^{[l]})$$

$$x = a^{[0]}$$

## Back propagation

- in put  $d_a^{[l]}$
  - out put  $d_a^{[l-1]}, d_w^{[l]}, d_b^{[l]}$
- $$\rightarrow d_z^{[l]} = d_a^{[l]} * g'^{[l]}(z^{[l]})$$
- $$\rightarrow d_w^{[l]} = d_z^{[l]} \cdot a^{[l-1].T}$$
- $$\rightarrow d_b^{[l]} = d_z^{[l]}$$
- $$\hookrightarrow d_a^{[l-1]} = w^{[l].T} \cdot d_z^{[l]}$$

## vectorized form

$$d_z^{[l]} = d_A^{[l]} * g'^{[l]}(z^{[l]})$$

$$d_w^{[l]} = (d_z^{[l]} \cdot A^{[l-1].T}) / m$$

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

$$d_A^{[l-1]} = w^{[l].T} \cdot d_z^{[l]}$$