

Deep Learning

Notes : Perceptron , Activation Function , Gradient Descent , backpropagation , loss Function , Metrics , Optimizer , Learning Rate , CNN , RNN , Structured Data , Unstructured Data , Research Papers , History

COURSE 1:
NN AND DEEPMLEARNING

Notes

→ ReLU : Rectified Linear Unit

LOGISTIC REGRESSION AS AN NN

1 / Notation

$$(x, y) \quad x \in \mathbb{R}^{n_x}, y \in \{0; 1\}$$

Training exp : $m = \{(x_1, y_1), \dots, (x_m, y_m)\}$

$$X = \begin{bmatrix} \vdots & & \\ x_1, \dots, x_m & & \\ \vdots & & \\ \xleftarrow[m]{} & & \xrightarrow[m_{\text{test}}]{} \end{bmatrix} \quad x_n$$

$$X \in \mathbb{R}^{n_x \times m}$$

$$X.\text{shape} = (n_x, m)$$

$$y = [y_1, \dots, y_m]$$

$$Y \in \mathbb{R}^{1 \times m}$$

$$Y.\text{shape} = (1, m)$$

Logistic Reg

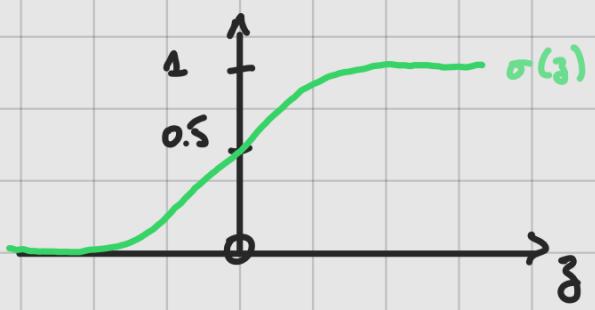
Given x , want $\hat{y} = P(y=1 | x)$

$$x \in \mathbb{R}^{n_x}$$

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

$$\text{output : } \hat{y} = \sigma(\underbrace{w^\top x + b}_{\text{Reg Lin}})$$

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



In some convention Define new extra feature :

$$x_0 = 1 \quad x \in \mathbb{R}^{n_x + 1}$$

$$\hat{y} = \sigma(\theta^T x)$$

avec

$$\theta = \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_{n_x} \end{bmatrix} \left\{ \begin{array}{l} b \\ w \end{array} \right\}$$

=> more easier to implement in NN

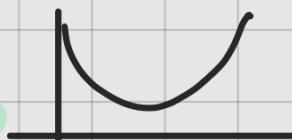
Logistic Regression Cost Function

$$\mathcal{L}(\hat{y}, y) = \frac{1}{3} (\hat{y} - y)^2 \text{ (MSE)} \rightarrow \text{Quadratic } \mathcal{W}$$

! Non convex

Want convex function (for Gradient descent) LOSS F

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



If $y = 1$: $\mathcal{L}(\hat{y}, 1) = -\log \hat{y}$ (want $\log \hat{y}$ large $\rightarrow \hat{y}$ large)

$y = 0$: $\mathcal{L}(\hat{y}, 0) = -\log(1-\hat{y})$ want \hat{y} small

Cost F: $J(\omega, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}_i, y_i)$

Gradient Descent



$$J(w)$$

Repeat

$$w := w - \alpha \frac{d J(w)}{d w}$$

↑
actualize

e_r

In code $d w =$

$$\hookrightarrow w := w - \alpha d w$$

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

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

In code $d w, d b$

Computation Graph

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

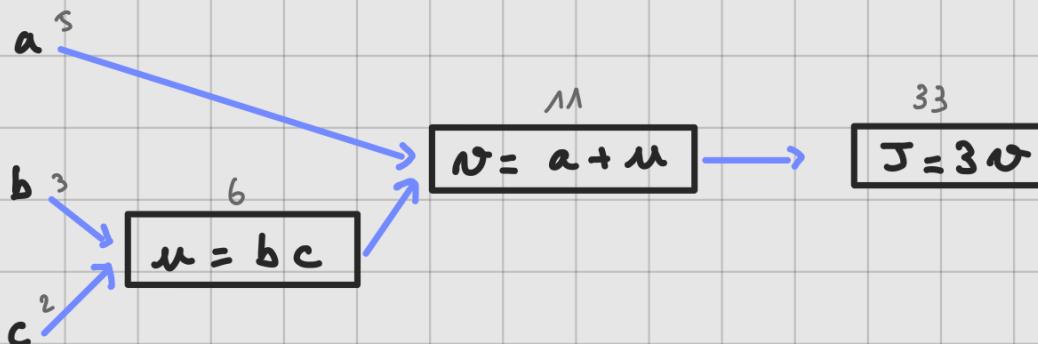
$\underbrace{}$
 u

$\underbrace{}$
 v

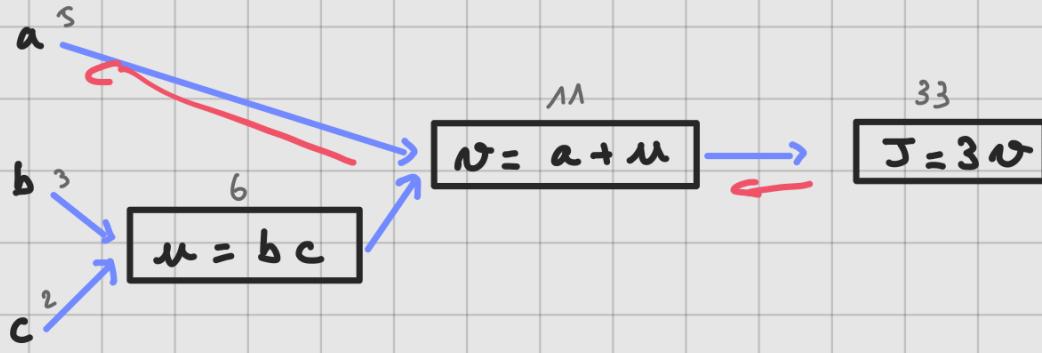
$$u = bc$$

$$v = a + u$$

$$J = 3v$$

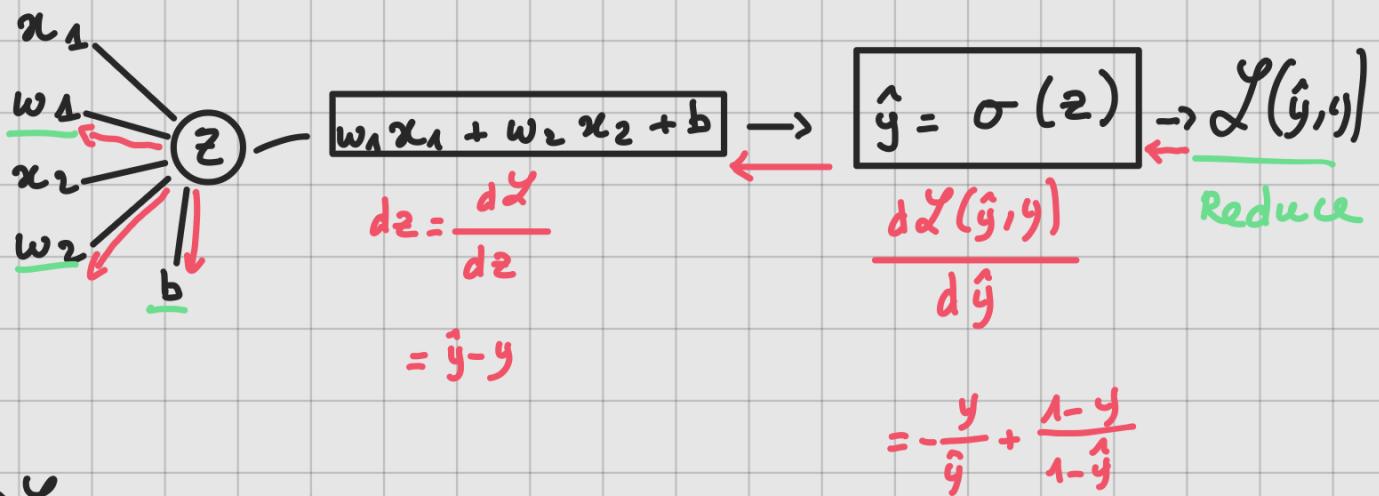


$$\frac{dJ}{dv} = ? = 3$$



Logistic Regression Gradient Descent

with 2 var : x_1, x_2



$$\frac{\partial L}{\partial w_1} = "dw_1" = x_1 dz$$

$$"dw_2" = x_2 dz$$

$$"db" = dz$$

$$\Rightarrow w_1 := w_1 - \alpha dw_1$$

$$w_2 := w_2 - \alpha dw_2$$

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

Logistic Regression on m examples

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

$$\hat{y}_i = \sigma(z_i) = \sigma(w^T x_i + b)$$

$$\frac{\partial}{\partial w_1} J(w, b) = \frac{1}{m} \sum_{i=1}^m \underbrace{\frac{\partial}{\partial w_i} \mathcal{L}(\hat{y}_i, y_i)}_{d w_a^{(i)}}$$

h algo:

$$J = 0, dw_1 = 0, dw_2 = 0, db = 0$$

For $i = 1$ to m :

$$z_i = w^T x_i + b$$

$$\hat{y}_i = \sigma(z_i)$$

$$J = - \left[y_i \log \hat{y}_i + (1-y_i) \log (1-\hat{y}_i) \right]$$

$$dz_i = \hat{y}_i - y_i$$

$$dw_a = x_a^{(i)} dz_i \quad \left. \right\} n=2$$

$$db := dz_i$$

$$\sum_j m_j = m ; \quad dw_1 / m_j = m ; \quad dw_2 / m_j = m ; \quad db / m_j = m$$

$$\text{update} \rightarrow w_1 := w_1 - \alpha dw_1$$

$$w_2 := w_2 - \alpha dw_2$$

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

Weakness: 2 for loop Δ

↳ have to vectorize

Vectorization

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

$$w = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \quad x = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \quad w \in \mathbb{R}^{n_x} \quad x \in \mathbb{R}^{n_x}$$

non vectorize:

$$z = 0$$

for i in range(n_x):

$$z += w[i] * x[i]$$

$$z += b$$

Vec:

$$z = np.\dot{w}(w, x) + b$$

much faster

Vectorizing Logistic Regression

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

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

$$\hat{y}^{(1)} = \sigma(z^{(1)})$$

$$\hat{y}^{(2)} = w^T x^{(2)} + b \dots$$

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

$$X \in \mathbb{R}^{n_x \times m}$$

shape = (n_x, m)

$$Z = [z^{(1)}, \dots, z^{(m)}] = w^T X + \underbrace{[b, b, \dots, b]}_{1 \times m}$$

$$= [w^T x^{(1)} + b, \dots, w^T x^{(m)} + b]$$

$$\hookrightarrow \underline{\text{np. dot}(w^T X) + b}$$

$$\hat{Y} = [\hat{y}^{(1)}, \dots, \hat{y}^{(m)}]$$

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

$$dZ = \hat{Y} - Y$$

$$dw = 0$$

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

$$dw += \dots$$

...

$$dw / \equiv m$$

$$db = 0$$

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

:

$$db / \equiv m$$

} For Loop

$$\rightarrow db = \frac{1}{m} \sum_{i=1}^m dz^{(i)}$$

$$= \frac{1}{m} \text{np.sum}(dz)$$

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

$$= \frac{1}{m} \begin{bmatrix} \vdots & \vdots \\ x^{(1)} & \cdots x^{(m)} \\ \vdots & \vdots \end{bmatrix} \begin{bmatrix} dz^{(1)} \\ \vdots \\ dz^{(m)} \end{bmatrix}$$

$$= \frac{1}{m} [x^{(1)} dz^{(1)}, \dots, x^{(m)} dz^{(m)}]$$

h algo:

$$J=0, dw_1=0, dw_2=0, db=0$$

For $i=1$ to m :

$$z_i = w^T x_i + b$$

$$\hat{y}_i = \sigma(z_i)$$

$$J += -[y_i \log \hat{y}_i + (1-y_i) \log (1-\hat{y}_i)]$$

$$dz_i = \hat{y}_i - y_i$$

$$\left. \begin{aligned} dw_1 &+= x_i^{(i)} dz_i \\ dw_2 &+= x_i^{(i)} dz_i \end{aligned} \right\}_{i=1}^m$$

$$dw += x^{(i)} * dz^{(i)}$$

$$db += dz_i$$

$$J/m ; dw_1 / m ; dw_2 / m ; db / m$$

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

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

$$\hat{Y} = \sigma(Z)$$

$$dZ = \hat{Y} - Y$$

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

$$db = \frac{1}{m} \text{np.sum}(dz)$$

$$w := w - \alpha dw$$

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

Broadcasting in Python

$$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + 100 \times \begin{bmatrix} 100 \\ 100 \\ 100 \\ 100 \end{bmatrix} = \begin{bmatrix} 101 \\ 102 \\ 103 \\ 104 \end{bmatrix}$$

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

$$+ \begin{bmatrix} 100 \\ 200 \end{bmatrix}_{(2,2)} \begin{bmatrix} 100 & 100 \\ 200 & 200 \end{bmatrix}_{(2,2)} = \begin{bmatrix} 101 & 102 & 103 \\ 204 & 205 & 206 \end{bmatrix}$$

$r + s$

$$(u \times v) + (w \times x)$$

$$(a+b)(a-b) + (b+c)(b-c)$$

$$a^2 - ab + ab - b^2 + b^2 - c^2$$

Algo

Forward and Backward Propagation

Inputs :

- w : weights
- b : bias
- X, Y

Outputs :

- Cost
- Grads

$$A = \sigma(w^T X + b)$$

forward

$$J = -\frac{1}{m} \sum_{i=1}^m (y^{(i)} \log(a^{(i)}) + (1-y^{(i)}) \log(1-a^{(i)}))$$

En numpy : $A = \text{sigmoid}(\text{np. dot}(w.T, X) + b)$

$$J = (1/m)^* \text{np. sum}(Y * \text{np. log}(A) + (1-Y) * \text{np. log}(1-A))$$

$$\frac{\partial J}{\partial w} = \frac{1}{m} \times (A - Y)^T$$

Backward

$$dw = \frac{1}{m} * \text{np. dot}(X, (A - Y).T)$$

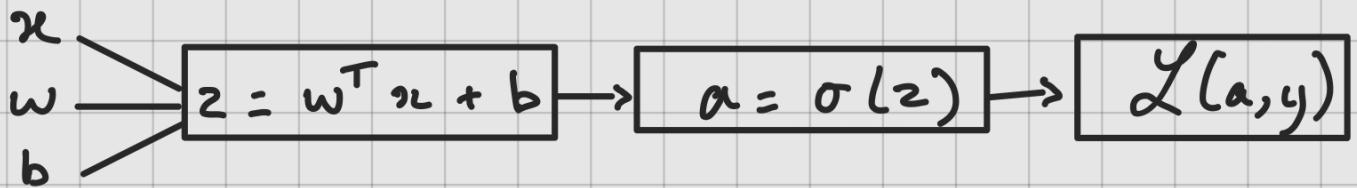
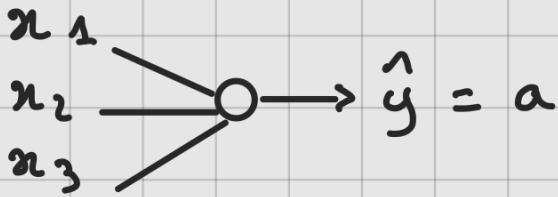
$$\frac{\partial J}{\partial b} = \frac{1}{m} \sum_{i=1}^m (a^{(i)} - y^{(i)})$$

$$db = \frac{1}{m} * \text{np. sum}(A - Y)$$

Return : Cost(J), Grads (dw, db)

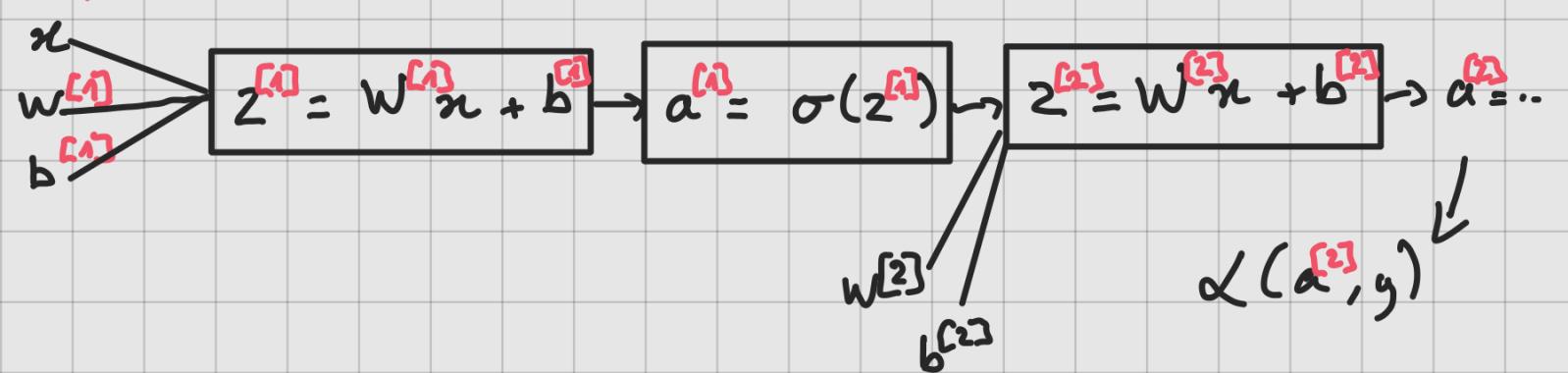
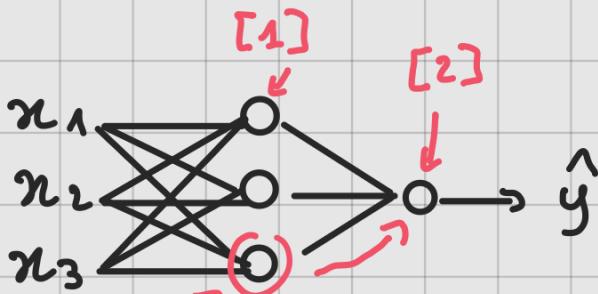
NEURAL NETWORK

Logistic :



Neural Network :

$[i]$ → layers
 (i) → examples

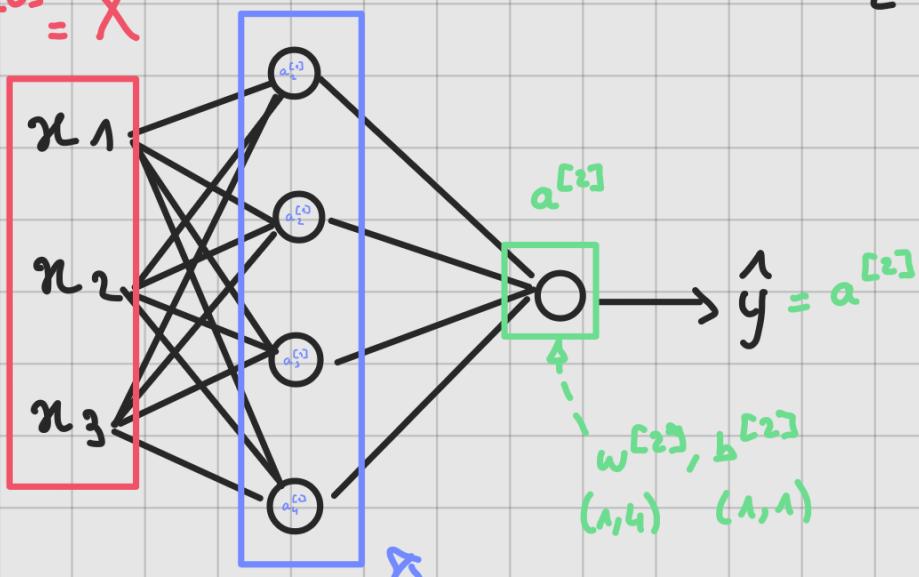


Notation :

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

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

2 layers NN (\emptyset Input)



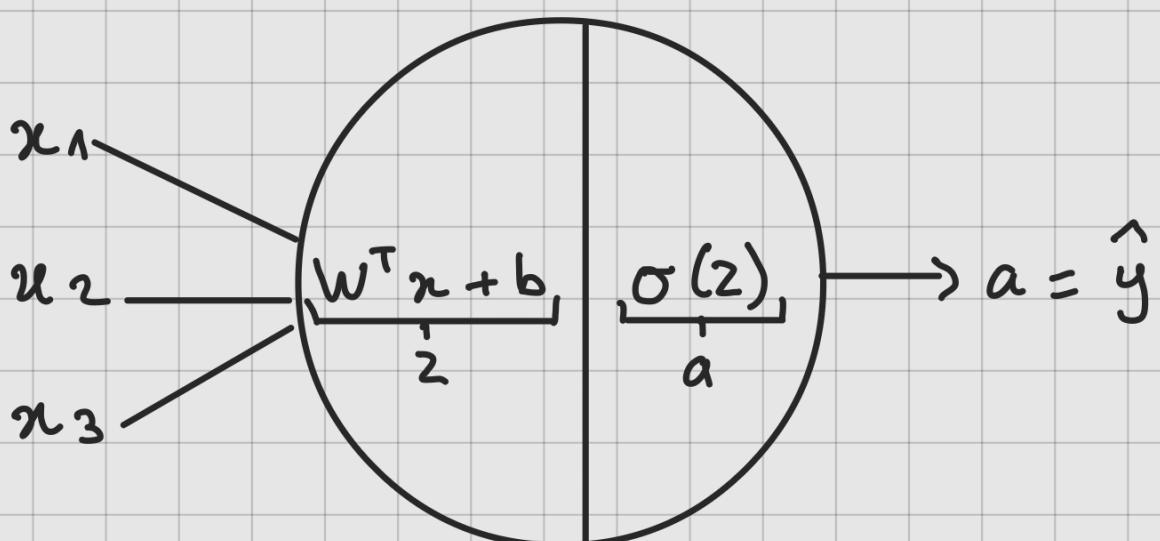
input
layer

hidden
layer

Output
layer

$$w^{[1]}, b^{[1]} \\ (4,3) \quad (4,1)$$

One nodes :



first N
 \downarrow
 fourth

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

$$\vdots$$

$$z_4^{[1]} = w_4^{[1]T} x + b_4^{[1]}, \quad a_4^{[1]} = \sigma(z_4^{[1]})$$

$$\begin{bmatrix}
 -w_1^{[1]T} \\
 -w_2^{[1]T} \\
 -w_3^{[1]T} \\
 -w_4^{[1]T}
 \end{bmatrix}_{(4,3)} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}_{(4,1)} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix}_{(4,1)} = \begin{bmatrix}
 w_1^{[1]T} x + b_1^{[1]} \\
 w_2^{[1]T} x + b_2^{[1]} \\
 w_3^{[1]T} x + b_3^{[1]} \\
 w_4^{[1]T} x + b_4^{[1]}
 \end{bmatrix}_{(4,1)}$$

$$z^{[1]} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}$$

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

Given an input x :

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

$(4,1) \quad (4,3) \quad (3,1) \quad (4,1)$

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

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

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

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

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

$(1,1) \quad (1,1)$

Vectorizing

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

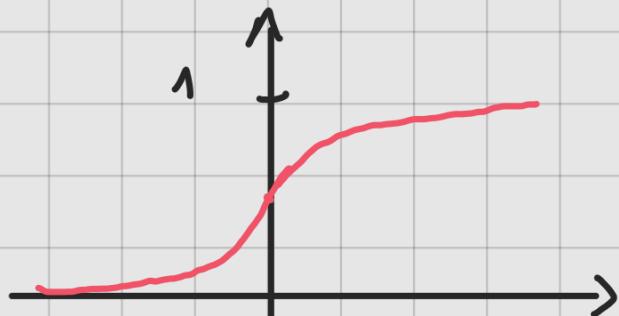
For $i = 1$ to m

$$z^{[1]}(i) = \dots$$
$$a^{[1]}(i) = \dots$$
$$a^{[2]}(i) = \dots$$

$x^{(i)} \longrightarrow a^{[2]}(i)$

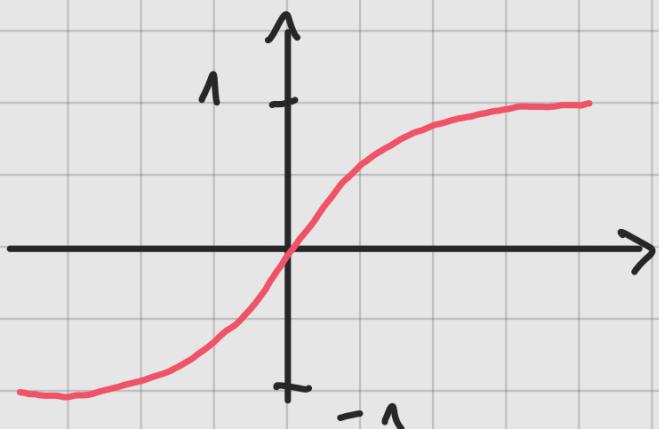
Activation Function

Sigmoid (used in output layer binary class)



$$a = \frac{1}{1 + e^{-x}}$$

Tanh : (almost always better) (except output layer)



$$a = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

ReLU :



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

↑ Leaky ReLU $a = \max(0.01z, z)$

Derivates

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

tanh

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

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

sigmoid

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

ReLU

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

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

$z \neq 0$

Leaky ReLU

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

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

Gradient Descent

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

$\underbrace{n_x = n^{[0]}}_{\text{nb input}}$, $\underbrace{n^{[1]}}_{\text{nb hidden}}$, $\underbrace{n^{[2]}}_{\text{output}} = 1$

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

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

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

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

Cost function :

$$J(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}, y)$$

Repeat :

Compute pred $(\hat{y}^{(1)}, \dots, \hat{y}^{(m)})$

$$dW^{[1]} = \frac{\partial J}{\partial W^{[1]}} , db^{[1]} = \frac{\partial J}{\partial b^{[1]}} , \dots$$

$$W^{[1]} := W^{[1]} - \alpha dW^{[1]}$$

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

:

Forward Propagation:

$$Z^{[1]} = W^{[1]} A^{[0]} + b^{[1]}$$

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

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

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

Back Propagation:

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

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

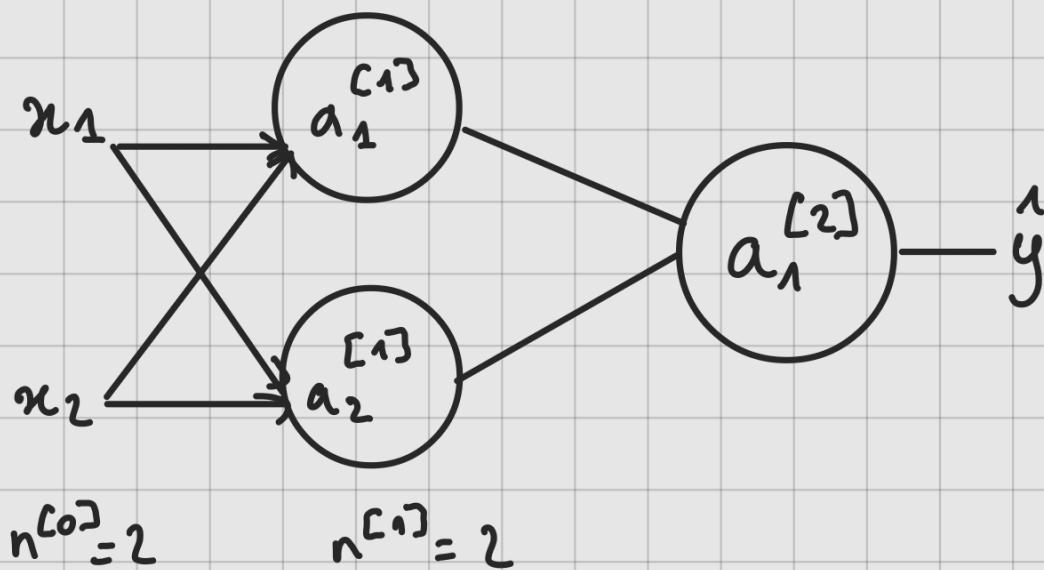
$$db^{[2]} = \frac{1}{m} \text{np.sum}(dZ^{[2]}, \text{axis}=1, \text{keepdims=True})$$

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

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

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

Random Initialization:



Init with 0 : $W^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$ $b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

$$\Rightarrow a_1^{[1]} = a_2^{[1]}, \quad d z_1^{[1]} = d z_2^{[1]}$$

Symmetric \rightarrow do exactly the same thing

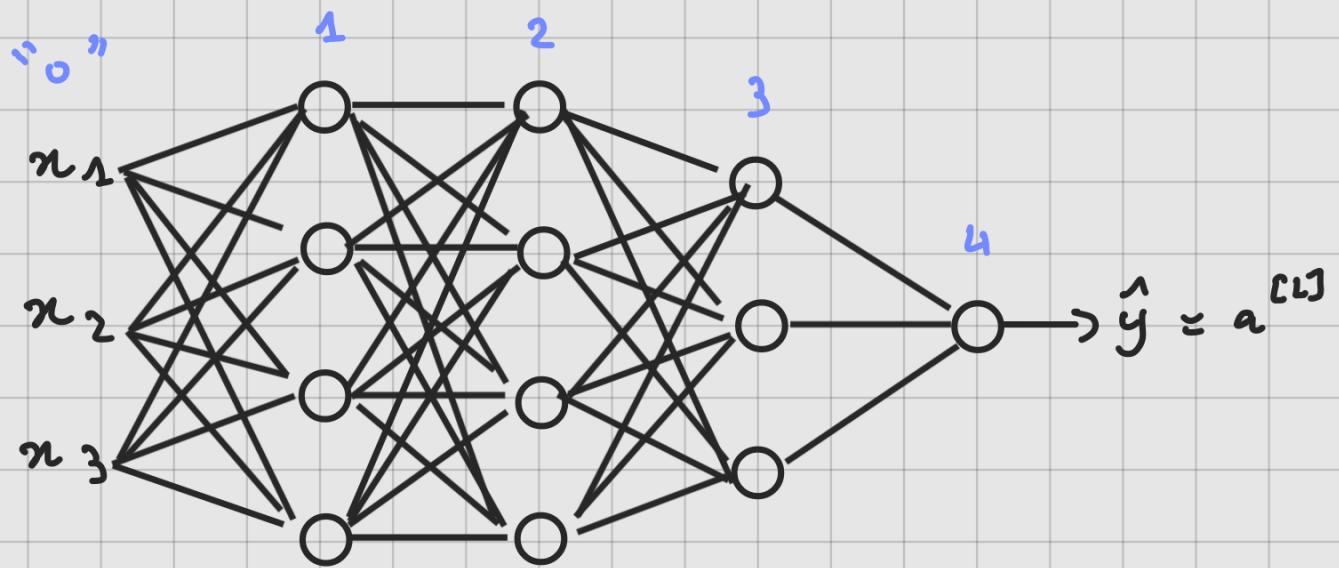
Random Init : $W^{[1]} = \text{np.random.randn}(2,2) * 0.01$

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

to have small values
 $\Rightarrow 2$
 GD faster

Deep NN

4 layer NN:



$L = 4$ (nb of layers)

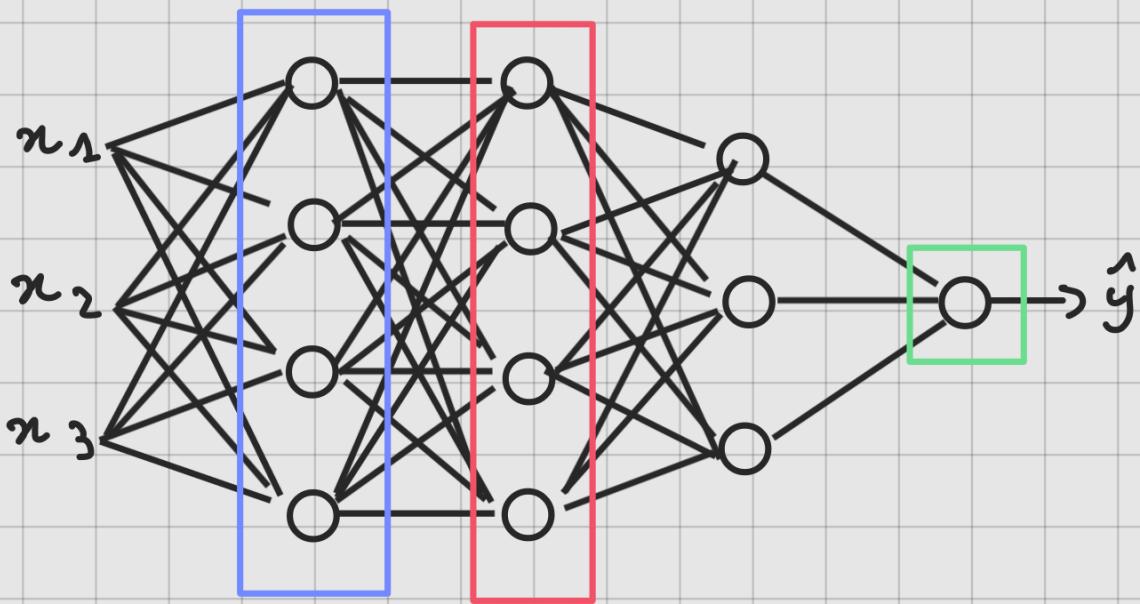
$n^{[l]}$ = nb of units in layer l

$$n^{[1]} = 5, \quad n^{[2]} = 5, \quad n^{[3]} = 3, \quad n^{[4]} = n^{[L]} = 1$$

$$n^{[0]} = n_x = 3 \quad x = a^{[0]}$$

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

Forward Propagation :



One output X :

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

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

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

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

...

$$Z^{[4]} = W^{[4]} a^{[3]} + b^{[4]}$$

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

$$Z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]}$$

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

Vectorize :

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

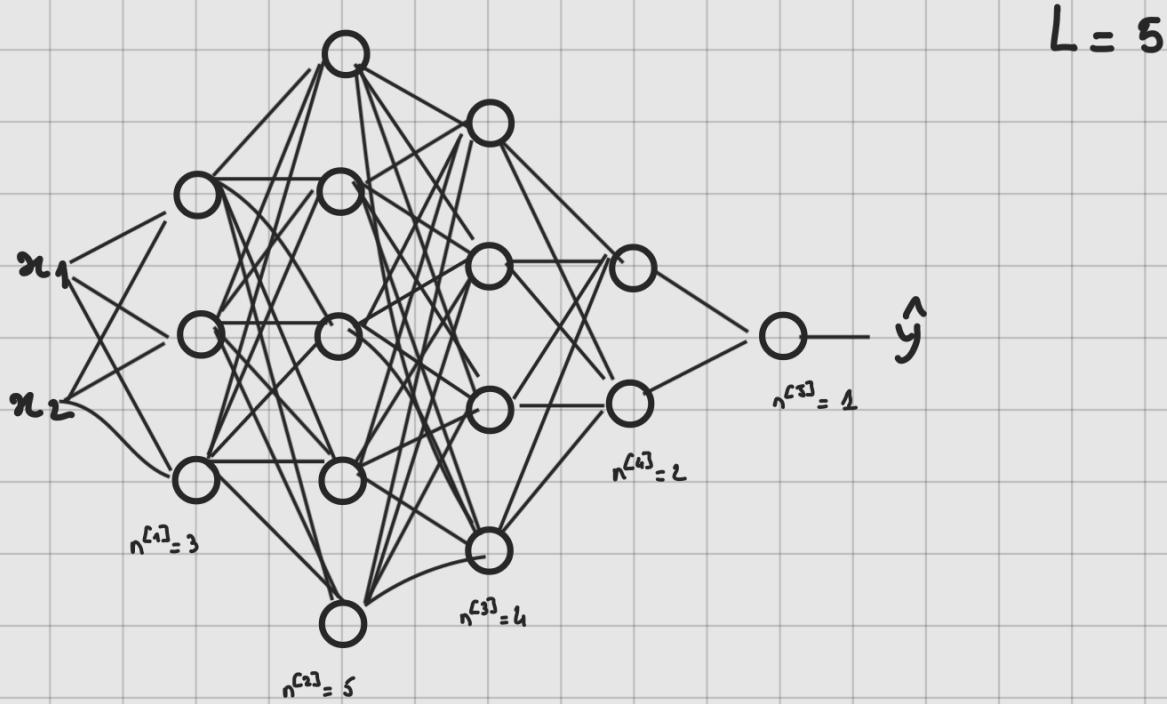
$$\hat{y} = g^{[4]}(Z^{[4]}) = A^{[4]}$$

...

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

↳ loop 1 to 4 (layers) for loop

Getting Matrix Dimensions Right



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

$$w^{[l]} = (n^{[l]}, n^{[l-1]})$$

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

$$(n^{[1]}, 1) \quad (n^{[2]}, n^{[1]}) \quad (n^{[3]}, 1)$$

$$\begin{bmatrix} \cdot \\ \vdots \\ \cdot \end{bmatrix} = \begin{bmatrix} \cdot & \cdot \\ \vdots & \vdots \\ \cdot & \cdot \end{bmatrix} \quad [\cdot]$$

$$a^{[l]}, z^{[l]}, b^{[l]} = (n^{[l]}, 1)$$

Vectorize

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

$\uparrow \quad \uparrow \quad \uparrow \quad \underbrace{\quad \quad \quad \quad}_{a}$
 $(n^{[1]}, m) \quad (n^{[1]}, n^{[0]}) \quad (n^{[0]}, m) \quad (n^{[1]}, 1)$
 $\xrightarrow{\text{Python broadcast}} (n^{[1]}, m)$

$$Z^{[e]}, A^{[e]} = (n^{[e]}, m)$$

Why deep representations?

take a look at Circuit theory and deep Learning

- There are functions you can compute with a "small" L-layer NN that shallower networks require exponentially more hidden units to compute.

Building Blocks of DNN:

Layer l : $W^{[l]}, b^{[l]}$

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

$$Z^{[L]} = W^{[L]} a^{[L-1]} + b^{[L]}$$

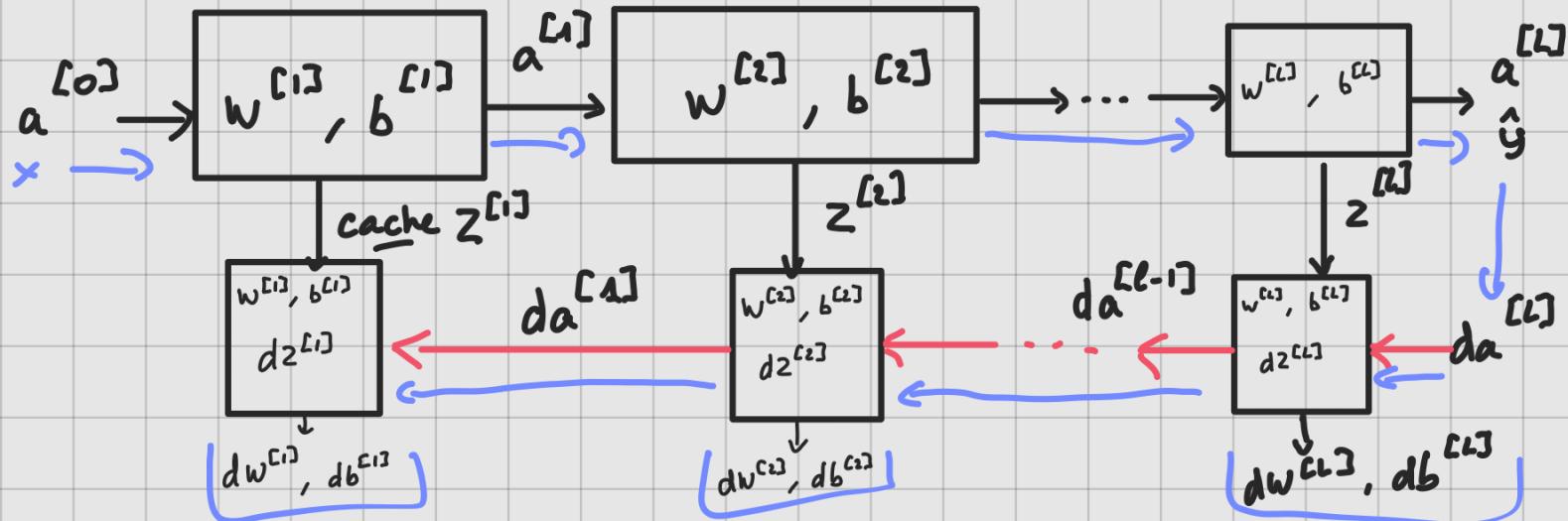
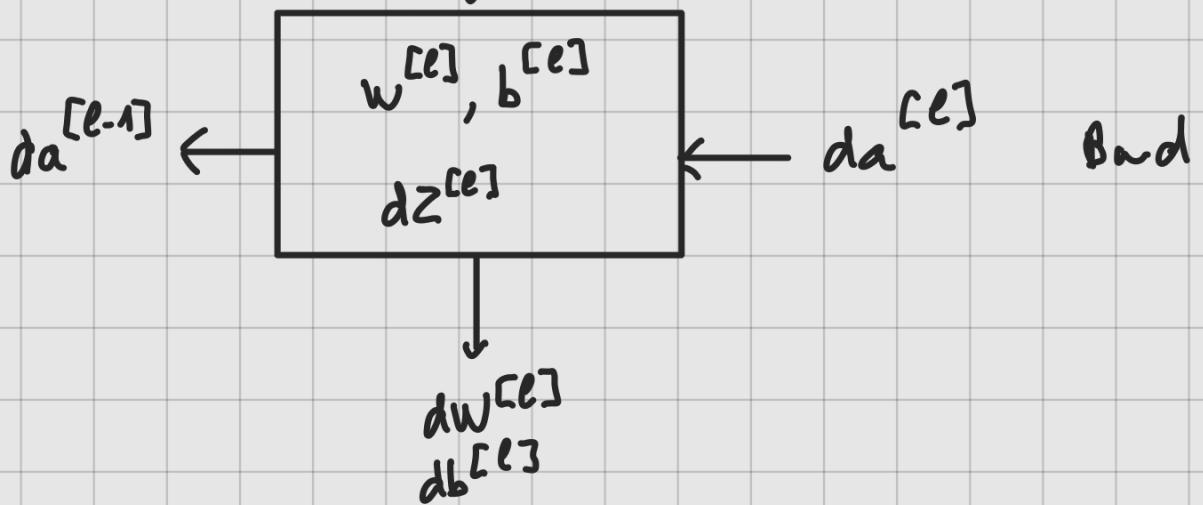
cache

$$Z^{[L]}$$

$$a^{[L]} = g^{[L]}(Z^{[L]})$$

Backward: Input $\frac{da^{[L]}}{\text{cache}(Z^{[L]})}$, output $\frac{da^{[L-1]}}{\frac{dW^{[L]}}{db^{[L]}}}$

Layer ℓ : $a^{[\ell-1]} \rightarrow [W^{[\ell]}, b^{[\ell]}] \rightarrow a^{[\ell]}$ Fund



$$w^{[e]} := w^{[e]} - \alpha dW^{[e]}$$

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

Forward Propagation:

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

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

$w^{[e]}, b^{[e]}$

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

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

vectorize

$$Z^{[l]} = W^{[l]} A^{[l-1]} + b^{[l]}$$

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

Backward Propagation:

Input: $da^{[l]}$

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

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

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

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

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

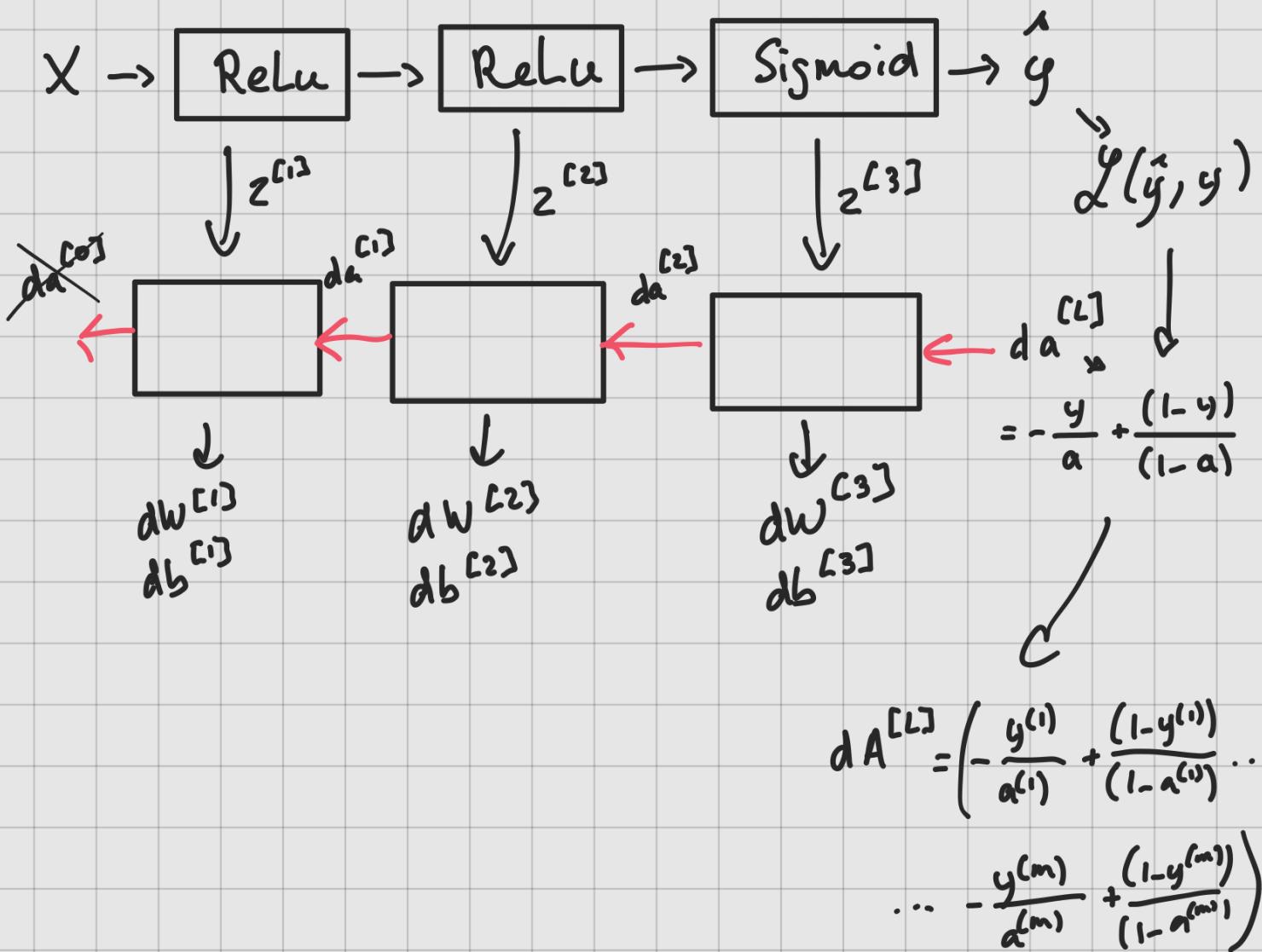
$$dz^{[l]} = dA^{[l]} * g^{[l]'}(Z^{[l]}) \quad \text{Vec}$$

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

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

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

Summarize :



Parameters vs Hyperparameters :

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

Hyperparameters : - Learning rate α ~

- Nb of iterations

- Nb of Hidden Layers L

- Nb of Hidden units $n^{[1]}, n^{[2]}, \dots$

- Choice of Activat^θ func

Other hyperparams : Momentum, minibatch size, regularization, ...