

# Computational graphs and backpropagation

software represent our neural net equations  
as graph

source node: inputs

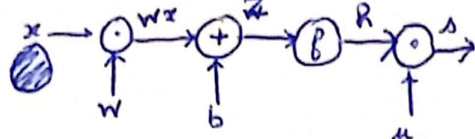
interior nodes: operations

$$s = u^T R$$

$$R = f(z)$$

$$z = wx + b$$

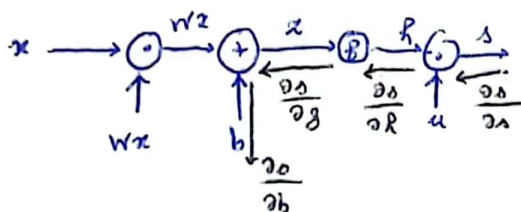
$$x: \text{input}$$



this is forward propagation

we want also to send back gradients to update the parameters of the model

⇓  
backpropagation



1 - f - propag:

visit nodes in topological sort order

- compute value of node given predecessors

2 - back-propag

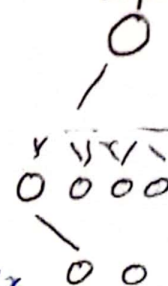
- initialize output gradient = 1

- visit nodes in reverse order

compute gradient of each node

using gradient wrt successors

$z$ : single scalar output

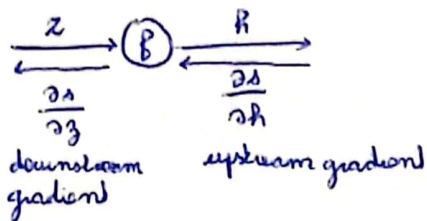


$\{y_1, y_2, \dots, y_n\}$ : successors of  $x$

$$\frac{\partial z}{\partial x} = \sum_{i=1}^n \frac{\partial z}{\partial y_i} \cdot \frac{\partial y_i}{\partial x}$$

backpropagation: single node

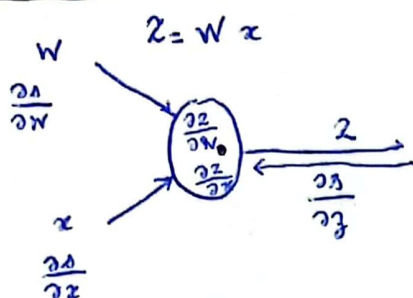
1/ input  $R = f(z)$



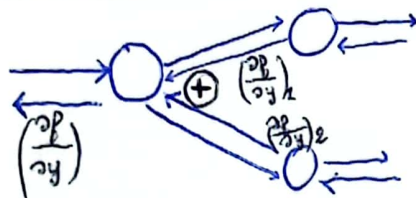
⇒ we need to calculate  $\frac{\partial R}{\partial z}$

$$\frac{\partial L}{\partial z} = \frac{\partial L}{\partial R} \cdot \frac{\partial R}{\partial z}$$

2/ multiple inputs



3/ multiple output branches



$$\left(\frac{\partial L}{\partial y}\right) = \left(\frac{\partial L}{\partial y}\right)_1 + \left(\frac{\partial L}{\partial y}\right)_2$$

⚠ incorrect way to compute backpropagation is to compute gradients one by one

⇒ we need to compute all gradients at once

# Math for neural networks

math 51 or Engr 408: matrix calculus

$$z = 3y^2 \quad \frac{dz}{dx} = \frac{dz}{dy} \cdot \frac{dy}{dx} = 6x$$

$$y = x^2$$

$$R = f(z) \Rightarrow \frac{\partial R}{\partial x} = \frac{\partial R}{\partial z} \cdot \frac{\partial z}{\partial x}$$

$$z = wx + b$$

$$R = f(z) \Rightarrow \frac{\partial R}{\partial z}$$

$$\text{if } z, R \in \mathbb{R}^n \Rightarrow \frac{\partial R}{\partial z} \in \mathbb{R}_{n,n} (\mathbb{R})$$

$$\left( \frac{\partial R}{\partial z} \right)_{ij} = \frac{\partial R_i}{\partial z_j}$$

1/ element wise transformations

like logistic transform

$$R = f(z) \Rightarrow R_i = f(z_i)$$

$$\left( \frac{\partial R}{\partial z} \right)_{ij} = \frac{\partial R_i}{\partial z_j} = \frac{\partial}{\partial z_j} f(z_i)$$

$$= \begin{cases} f'(z_i) & \text{if } i=j \\ 0 & \text{else} \end{cases}$$

$$\frac{\partial R}{\partial z} = \begin{pmatrix} f'(z_1) & & (0) \\ & \ddots & \\ (0) & & f'(z_n) \end{pmatrix}$$

$$\left. \begin{array}{l} X \in \mathbb{R}^m \\ W \in \mathbb{R}_{n,m} (\mathbb{R}) \\ b \in \mathbb{R}^n \end{array} \right\} z \in \mathbb{R}^n$$

$$\left. \begin{array}{l} \text{element wise transformation} \end{array} \right\} R \in \mathbb{R}^n \quad R = f(z)$$

$$u \in \mathbb{R}^n \Rightarrow v \in \mathbb{R}$$

$$2/ \frac{\partial}{\partial x} (wx + b) = W$$

$$3/ \frac{\partial}{\partial b} (wx + b) = I$$

$$4/ \frac{\partial}{\partial u} (u^T R) = R^T \quad \frac{\partial}{\partial R} (u^T R) = u^T$$

$$v = u^T R$$

$$R = f(Wx + b) = f(z) \quad z = wx + b$$

$$x = \text{input}$$

$$\frac{\partial v}{\partial b} = \frac{\partial v}{\partial R} \cdot \frac{\partial R}{\partial z} \cdot \frac{\partial z}{\partial b}$$

$$= u^T \cdot \text{diag}(f'(z)) \cdot I$$

$$= u^T \circ \text{diag}(f'(z))$$

Radamand product

$$W \in \mathbb{R}^{n \times m} \Rightarrow \frac{\partial v}{\partial W} : 1 \text{ by } n \times m \text{ jacobian}$$

it is a huge row vector

$$\theta^{\text{new}} = \theta^{\text{old}} - \alpha \cdot \nabla_{\theta} J(\theta)$$

inconvenient to do so if we have a row vector

$\Rightarrow$  we apply the shape convention

$$1 \times n \times n/n \times n, n, m$$

~~Revisiting the shape convention~~

## shape conventions

$\frac{\partial \lambda}{\partial W} \in \mathbb{R}^{1, nm}(\mathbb{R})$  : row vector

$\Rightarrow$  we can't calculate  $\delta^{new} = \delta^{old} - \alpha \cdot \nabla_{\theta} J(\theta)$

$\Rightarrow$  we leave pure math and we use the shape convention

the shape of the gradient is the shape of the parameters

$\Rightarrow W \in \mathbb{R}^{n, m}(\mathbb{R}) \Rightarrow \frac{\partial \lambda}{\partial W} \in \mathbb{R}^{1, nm}(\mathbb{R})$

$$\frac{\partial \lambda}{\partial W} = \begin{bmatrix} \frac{\partial \lambda}{\partial w_{11}} & \dots & \frac{\partial \lambda}{\partial w_{1m}} \\ \vdots & \ddots & \vdots \\ \frac{\partial \lambda}{\partial w_{n1}} & \dots & \frac{\partial \lambda}{\partial w_{nm}} \end{bmatrix}$$

$$\frac{\partial \lambda}{\partial W} = \frac{\partial \lambda}{\partial R} \cdot \frac{\partial R}{\partial z} \cdot \frac{\partial z}{\partial W}$$

$$S = u^T \cdot \text{diag}(\phi'(z))$$

$1 \times n, nm$

$$= S \cdot \left( \frac{\partial z}{\partial W} \right) \in \mathbb{R}^{1, nm}(\mathbb{R})$$

analog

the answer is  $\left( \frac{\partial \lambda}{\partial W} \right)_{\text{convention}} = S^T x^T$  why the transpose?

$$\frac{\partial \lambda}{\partial W} = S \cdot \frac{\partial \lambda}{\partial W} (Wx + b)$$

for a single weight  $w_{ij}$

$$\frac{\partial z_i}{\partial w_{ij}} = \frac{\partial}{\partial w_{ij}} w_i x + b_i$$

$$= \frac{\partial}{\partial w_{ij}} \sum_{k=1}^d w_{ik} x_k = x_j$$

$\rightarrow$  we want  $\frac{\partial \lambda}{\partial W} \in \mathbb{R}^{1, nm}(\mathbb{R})$  and not  $\mathbb{R}^{1, nm}(\mathbb{R})$

$\Rightarrow S \in \mathbb{R}^{1, n}(\mathbb{R}) \quad x \in \mathbb{R}^{m, 1}(\mathbb{R})$

$\Rightarrow$  by transposing both we obtain the result

$\bullet \frac{\partial \lambda}{\partial b} = u^T \phi'(z)$  is a row vector

shape convention says that our gradient should have the same shape of  $b \Rightarrow$  a column vector

$$\left( \frac{\partial \lambda}{\partial b} \right)_{\text{convention}} = (u^T \phi'(z))^T$$

Jacobian form is useful to do calculus



to do so !

1/ use jacobian form as much as possible, reshape to follow the shape convention at the end (what we did for  $\frac{\partial \lambda}{\partial b}$ )

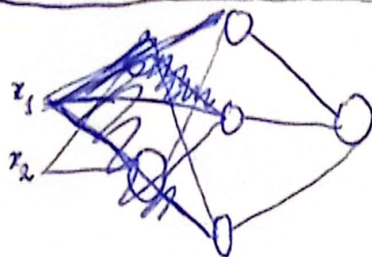
2/ always follow the shape convention

Look at dimensions to figure out when to transpose and/or rescale terms





## Reseau de neurones à 2 couches



$$X = \begin{bmatrix} x_1^{(1)} & \dots & x_1^{(m)} \\ x_2^{(1)} & \dots & x_2^{(m)} \end{bmatrix}$$

$m$ : nombre d'observations  
ici on a 2 variables;  $n^{[0]} = 2$

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

$$A^{[1]} = \frac{1}{1 + e^{-z^{[1]}}}$$

$$W \in \mathbb{R}^{n^{[1]}, n^{[0]}}$$

$$b \in \mathbb{R}^{n^{[1]}, 1}$$

$$z^{[1]}, b^{[1]} \in \mathbb{R}^{n^{[1]}, m}$$

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

$$A^{[2]} = \frac{1}{1 + e^{-z^{[2]}}}$$

fonction coût dans ce cas

$$\mathcal{L} = -\frac{1}{m} \sum y \log(A^{[2]}) + (1-y) \log(1 - A^{[2]})$$

## Backpropagation

$$\frac{\partial \mathcal{L}}{\partial W^{[2]}} = \frac{\partial \mathcal{L}}{\partial A^{[2]}} \cdot \frac{\partial A^{[2]}}{\partial z^{[2]}} \cdot \frac{\partial z^{[2]}}{\partial W^{[2]}}$$

$$\frac{\partial \mathcal{L}}{\partial b^{[2]}} = \frac{\partial \mathcal{L}}{\partial A^{[2]}} \cdot \frac{\partial A^{[2]}}{\partial z^{[2]}} \cdot \frac{\partial z^{[2]}}{\partial b^{[2]}}$$

$$\frac{\partial \mathcal{L}}{\partial W^{[1]}} = \frac{\partial \mathcal{L}}{\partial A^{[1]}} \cdot \frac{\partial A^{[1]}}{\partial z^{[1]}} \cdot \frac{\partial z^{[1]}}{\partial W^{[1]}}$$

$$\frac{\partial \mathcal{L}}{\partial b^{[1]}} = \frac{\partial \mathcal{L}}{\partial A^{[1]}} \cdot \frac{\partial A^{[1]}}{\partial z^{[1]}} \cdot \frac{\partial z^{[1]}}{\partial b^{[1]}}$$

$$\frac{\partial \mathcal{L}}{\partial W^{[1]}} = \frac{\partial \mathcal{L}}{\partial A^{[2]}} \cdot \frac{\partial A^{[2]}}{\partial z^{[2]}} \cdot \frac{\partial z^{[2]}}{\partial A^{[1]}} \cdot \frac{\partial A^{[1]}}{\partial z^{[1]}} \cdot \frac{\partial z^{[1]}}{\partial W^{[1]}}$$

$$\frac{\partial \mathcal{L}}{\partial b^{[1]}} = \frac{\partial \mathcal{L}}{\partial A^{[2]}} \cdot \frac{\partial A^{[2]}}{\partial z^{[2]}} \cdot \frac{\partial z^{[2]}}{\partial A^{[1]}} \cdot \frac{\partial A^{[1]}}{\partial z^{[1]}} \cdot \frac{\partial z^{[1]}}{\partial b^{[1]}}$$

$$d_{\mathcal{L}} = \frac{\partial \mathcal{L}}{\partial A^{[2]}} \cdot \frac{\partial A^{[2]}}{\partial z^{[2]}}$$

$$= \frac{1}{m} \sum \left( -\frac{y}{A^{[2]}} + \frac{(1-y)}{1 - A^{[2]}} \right) \times A^{[2]} (1 - A^{[2]})$$

$$= \frac{1}{m} \sum -y(1 - A^{[2]}) + \frac{1}{m} \sum (1-y) A^{[2]}$$

$$= \frac{1}{m} \sum -y + y A^{[2]} + A^{[2]} - y A^{[2]}$$

$$d_{\mathcal{L}} = \frac{1}{m} \sum A^{[2]} - y$$

$\downarrow$   $\downarrow$   $\downarrow$   
 $(n^{[2]}, m)$   $(n^{[2]}, m)$   $(1, m)$

on fait un broadcasting pour  $y$

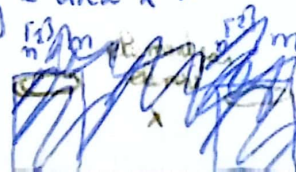
$$A^{[2]} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}_{n^{[2]} \times m} - \begin{bmatrix} y \\ y \\ y \end{bmatrix}_{1 \times m}$$

$\Rightarrow$  on crée une matrice  $(n^{[2]}, m)$  en répétant le

$$\text{vecteur } y$$

$\begin{pmatrix} y \\ y \\ y \end{pmatrix}_{(n^{[2]}, m)}$   $\begin{pmatrix} y \\ y \\ y \end{pmatrix}_{(n^{[2]}, m)}$   $\begin{pmatrix} y \\ y \\ y \end{pmatrix}_{(n^{[2]}, m)}$

$$\text{alors } \frac{\partial \mathcal{L}}{\partial W^{[2]}} = d_{\mathcal{L}} \times A^{[1]T}$$



## Gradient descent for NN

1. one hidden layer NN for binary classification

parameters:  $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}$   
 $n = n^{[0]}, n^{[1]}, n^{[2]} = 1$   
 (n<sup>[1]</sup>, n<sup>[0]</sup>) (n<sup>[2]</sup>, 1) (n<sup>[1]</sup>, n<sup>[2]</sup>)

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

### Gradient descent

Repeat {

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

$$dw^{[1]} = \frac{\partial J}{\partial w^{[1]}}$$

$$db^{[1]} = \frac{\partial J}{\partial b^{[1]}}$$

$$w^{[1]} = w^{[1]} - \alpha \cdot dw^{[1]}$$

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

### Random initialization

if we initialize all weights to zero

for any example  $a_1^{[1]} = a_1^{[2]}$

$$\Rightarrow dz_1^{[1]} = dz_1^{[2]} \quad (\text{symmetry breaking problem})$$



best initialization

$$w^{[1]} = \text{np.random.randn}(4, 3) * 0.01$$

$$b^{[1]} = \text{np.zeros}(4, 1) \quad (\text{no pb for } b)$$

formulas for computing derivatives

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

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

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

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

### Backpropagation

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

$$dw^{[2]} = \frac{1}{m} dz^{[2]} A^{[1]T}$$

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

keepdims = True

prevent numpy to return  
rank 1 arrays

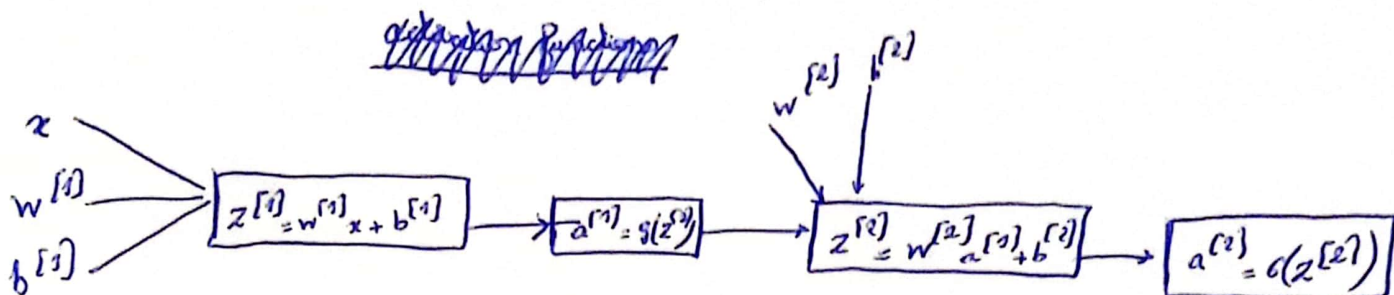
$$dz^{[1]} = w^{[2]T} dz^{[2]} * g^{[1]'}(z^{[1]})$$

element wise product

$$dw^{[1]} = \frac{1}{m} dz^{[1]} x^T$$

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

to have small numbers to avoid flat parts  
of the activation functions  $\Rightarrow$  faster computation



$$\mathcal{L}(a^{[2]}, y) = -y \log(a^{[2]}) + (1-y) \log(1-a^{[2]})$$

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

$$dz^{[2]} = da^{[2]} \cdot \frac{da}{dz^{[2]}}$$

$$= da^{[2]} \cdot g'(z) \quad \text{or } g(z) = c(z)$$

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

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

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

$$da^{[1]} = w^{[2]T} \cdot dz^{[2]}$$

$$dz^{[1]} = \cancel{w^{[2]T}} w^{[2]T} \cdot dz^{[2]} * g^{[1]'}(z^{[1]})$$

d

$$dw^{[1]} = dz^{[1]} a^{[0]T}$$

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

element  
wise product

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

$$z^{[i]} = w^{[i]}x + b^{[i]}$$

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

↓

vectorized

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

$$dw^{[2]} = \frac{1}{m} dz^{[2]} A^{[1]T}$$

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

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

$$dw^{[1]} = \frac{1}{m} dz^{[1]} A^{[0]T}$$

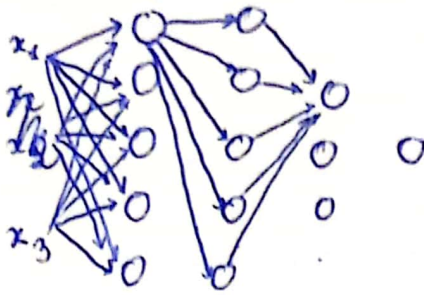
$$db^{[1]} = \frac{1}{m} \cdot \text{np.sum}(dz^{[1]}, \text{axis}=1, \text{keepdims}=\text{True})$$



## Deep L-layer NN

The number of hidden layers is a hyperparameter  
 The number of neurons in hidden layers is a hyperparameter

example: 4 layer neural network



$$L = \# \text{ layers} = 4$$

$$n^{[L]} = \# \text{ units in layer } L \in \{1, \dots, L\}$$

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

$$n^{[0]} = n = \# \text{ features}$$

$$a^{[L]} = \text{activation in layer } L$$

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

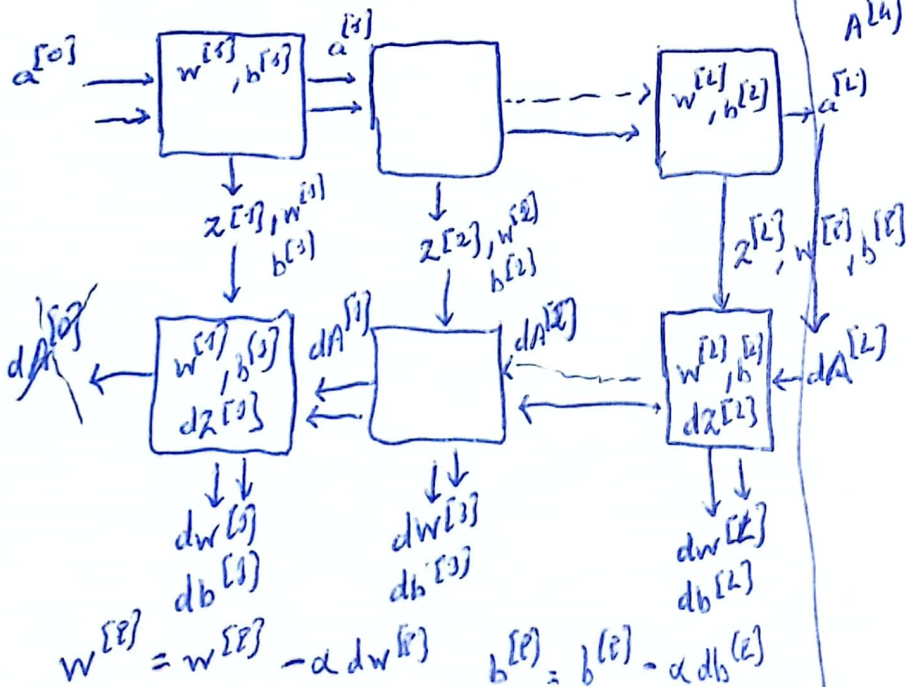
$$w^{[L]} = \text{weights for } z^{[L]}$$

$$b^{[L]}$$

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

$$a^{[L]} = \hat{y}$$

computation graph for forward and backward propagation (one iteration of gradient descent)



forward propagation for a single training example

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

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

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

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

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

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

$$\Rightarrow z^{[L]} = w^{[L]T} a^{[L-1]} + b^{[L]}$$

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

forward propagation for the whole training set

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

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

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

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

$$Z^{[3]} = w^{[3]} A^{[2]} + b^{[3]}$$

$$A^{[3]} = g^{[3]}(Z^{[3]})$$

$$Z^{[4]} = w^{[4]} A^{[3]} + b^{[4]}$$

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

On n'a pas le choix entre de faire un for loop over the L layers of the NN



## Forward propagation

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

output  $a^{[P]}$ , cache  $(z^{[P]}, w^{[P]}, b^{[P]})$

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

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

## backward

$$dz^{[P]} = dA^{[P]} \underset{\substack{\text{element wise} \\ \text{product}}}{*} g^{[P]'}(z^{[P]})$$

$$dA^{[P]} = \cancel{dA^{[P]}} \underset{\substack{\text{element wise} \\ \text{product}}}{*} \frac{1}{m} dz^{[P]} A^{[P-1]T}$$

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

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

# TD 6 : Réseaux de neurones

Def 1

$$g: \mathbb{R}^n \times \mathbb{R}^n \times \mathbb{R} \rightarrow \mathbb{R}$$

$$(x, w, b) \mapsto g(x, w, b) = \sigma\left(\sum_{i=1}^n w_i x_i + b\right)$$

$\sigma$  : fonction d'activation

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

Def 2 : apprentissage du neurone : estimation par les moindres carrés des paramètres du neurone

1)  $\left\{ \begin{array}{l} r: \mathbb{R}^n \times \mathbb{R} \rightarrow \mathbb{R}^K \\ (w, b) \end{array} \right. \longrightarrow \left( \begin{array}{l} y^1 - \sigma\left(\sum_{i=1}^n w_i x_i + b\right) \\ \vdots \\ y^K - \sigma\left(\sum_{i=1}^n w_i x_i + b\right) \end{array} \right)$

$n(w_i)$  → on a K points  
indice pas puissance  
indice pas puissance

$$\left\{ \begin{array}{l} \min f(w, b) = \frac{1}{2} \|r(\beta)\|^2 \\ (w, b) \in \mathbb{R}^n \times \mathbb{R} \end{array} \right.$$

2) ce n'est pas linéaire car  $\sigma$  est non linéaire

3)  $\sigma = \text{Id}$

$$\begin{aligned} r(w, b) &= \begin{pmatrix} y^1 - \sum_{i=1}^n w_i x_i + b \\ \vdots \\ y^K - \sum_{i=1}^n w_i x_i + b \end{pmatrix} \\ &= \begin{pmatrix} y^1 \\ \vdots \\ y^K \end{pmatrix} - \begin{pmatrix} (x^1)^T & 1 \\ \vdots & 1 \\ (x^K)^T & 1 \end{pmatrix} \begin{pmatrix} w \\ b \end{pmatrix} \end{aligned}$$

$$X = \begin{pmatrix} (x^1)^T & 1 \\ \vdots & 1 \\ (x^K)^T & 1 \end{pmatrix}$$

$$x^i = \begin{pmatrix} x_1^i \\ \vdots \\ x_n^i \end{pmatrix}$$

$$w = \begin{pmatrix} w_1 \\ \vdots \\ w_n \end{pmatrix}$$

$$1.4 / r(\beta) = \begin{pmatrix} y_0^1 - \sigma(\langle w, x^1 \rangle + b) \\ \vdots \\ y^K - \sigma(\langle w, x^K \rangle + b) \end{pmatrix}$$

$$\nabla_{\beta}(\beta) = J_{\beta}(\beta)^T \cdot r(\beta)$$

$$J_{\beta}(\beta) = \begin{pmatrix} -\sigma'(\langle w, x^1 \rangle + b) (x^1)^T & -\sigma'(\langle w, x^1 \rangle + b) \\ \vdots & \vdots \\ -\sigma'(\langle w, x^K \rangle + b) (x^K)^T & -\sigma'(\langle w, x^K \rangle + b) \end{pmatrix}$$

## Exercice 2 : réseau à plusieurs couches

2.1 / ~~nn~~  $L$  : nombre de couches = 4

$l: 1, \dots, L$  couches

$$n_1 = 2$$

$$n_2 = 5$$

$$n_3 = 3$$

$$n_4 = 1$$

$$n_0 = 4$$

~~W~~ et de dimension  $L$  (car  $n_0 = 4$ )

$$W^1 \in M_{(n_1, n_0)}(\mathbb{R})$$

$$b^1 \in \mathbb{R}^{n_1} = \begin{bmatrix} b_1^1 \\ b_2^1 \end{bmatrix}$$

$$W^2 \in M_{(n_2, n_1)}(\mathbb{R})$$

$$b^2 \in \mathbb{R}^{n_2} = 5$$

$$W^3 \in M_{(n_3, n_2)}(\mathbb{R})$$

$$b^3 \in \mathbb{R}^3$$

$$W^4 \in M_{(n_4, n_3)}(\mathbb{R})$$

$$b^4 \in \mathbb{R}^1$$

2.2 /  $a^l$  : sortie de la couche  $l$  : activation de la couche  $l$

→ on applique  $\sigma$  à chaque élément du vecteur

$$a^{l+1} = \sigma_{n_{l+1}}(W^{l+1} \cdot a^l + b^{l+1})$$

$$= \begin{pmatrix} \sigma \\ \vdots \\ \sigma \end{pmatrix}$$



$$\beta = \{ w^1, b^1; w^2, b^2; w^3, b^3; w^4, b^4 \}$$

$$\in \Pi_{f, n_0}(\mathbb{R}) \times \mathbb{R}^{n_1} \times \Pi_{g, n_1}(\mathbb{R}) \times \mathbb{R}^{n_2} \times \Pi_{h, n_2}(\mathbb{R}) \times \mathbb{R}^{n_3} \times \Pi_{i, n_3}(\mathbb{R}) \times \mathbb{R}^{n_4} = \mathcal{E}$$

~~ensemble~~

alors  $y \left( \underset{\substack{\mathbb{R}^k \\ \mathcal{E}}}{x}, \underset{\mathcal{E}}{\beta} \right) = a^4 = \vec{\sigma}_{n_4} \left( w^4 \vec{\sigma}_{n_3} \left( w^3 \vec{\sigma}_{n_2} \left( w^2 \vec{\sigma}_{n_1} (w^1 x + b^1) + b^2 \right) + b^3 \right) + b^4 \right)$

terminaison de tous les neurones

2.3 /  $r: \mathcal{E} \rightarrow \mathbb{R}^K$

$$r(\beta) = \begin{pmatrix} y^k - y(x^k, \beta) \\ \vdots \\ y^K - y(x^K, \beta) \end{pmatrix}$$

$$\begin{cases} \min_{\beta \in \mathcal{E}} J(\beta) = \frac{1}{2} \|r(\beta)\|^2 \end{cases}$$

2.4 / ~~der~~ dérivées partielles par rapport à  $\beta^4 = (w^4, b^4)$

$$\frac{\partial y}{\partial \beta^4} = \vec{\sigma}'_{n_4} \left( w^4 \vec{\sigma}_{n_3} \left( w^3 \vec{\sigma}_{n_2} \left( w^2 \vec{\sigma}_{n_1} (w^1 x + b^1) + b^2 \right) + b^3 \right) + b^4 \right) \odot \left[ \vec{\sigma}_{n_3} \left( w^3 \vec{\sigma}_{n_2} \left( w^2 \vec{\sigma}_{n_1} (w^1 x + b^1) + b^2 \right) + b^3 \right)^T, 1 \right]$$

produit terme à terme

$$\frac{\partial y}{\partial \beta^3} = \vec{\sigma}'_{n_4} \left( w^4 \vec{\sigma}_{n_3} \left( w^3 \vec{\sigma}_{n_2} \left( w^2 \vec{\sigma}_{n_1} (w^1 x + b^1) + b^2 \right) + b^3 \right) + b^4 \right) \odot \left[ \frac{\partial}{\partial \beta^3} \left( w^4 \vec{\sigma}_{n_3} \left( w^3 \vec{\sigma}_{n_2} \left( w^2 \vec{\sigma}_{n_1} (w^1 x + b^1) + b^2 \right) + b^3 \right) + b^4 \right) \right]$$

||

$$\left[ w^4 \left( \frac{\partial}{\partial \beta^3} \left( \vec{\sigma}_{n_3} \left( w^3 \vec{\sigma}_{n_2} \left( w^2 \vec{\sigma}_{n_1} (w^1 x + b^1) + b^2 \right) + b^3 \right) \right) \right) \right]$$