

# Neural Network

Deep learning par la pratique

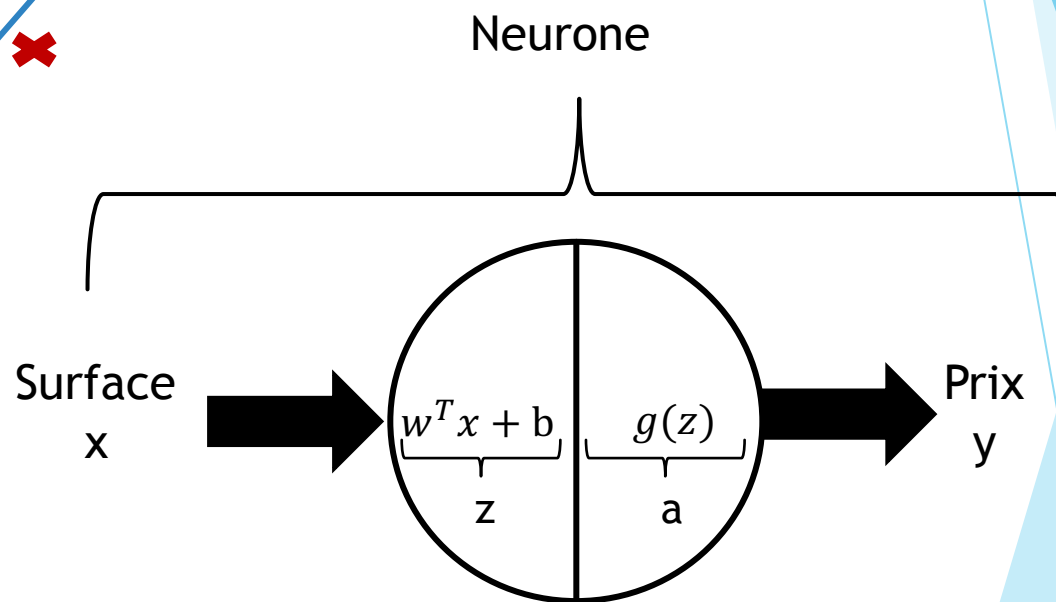
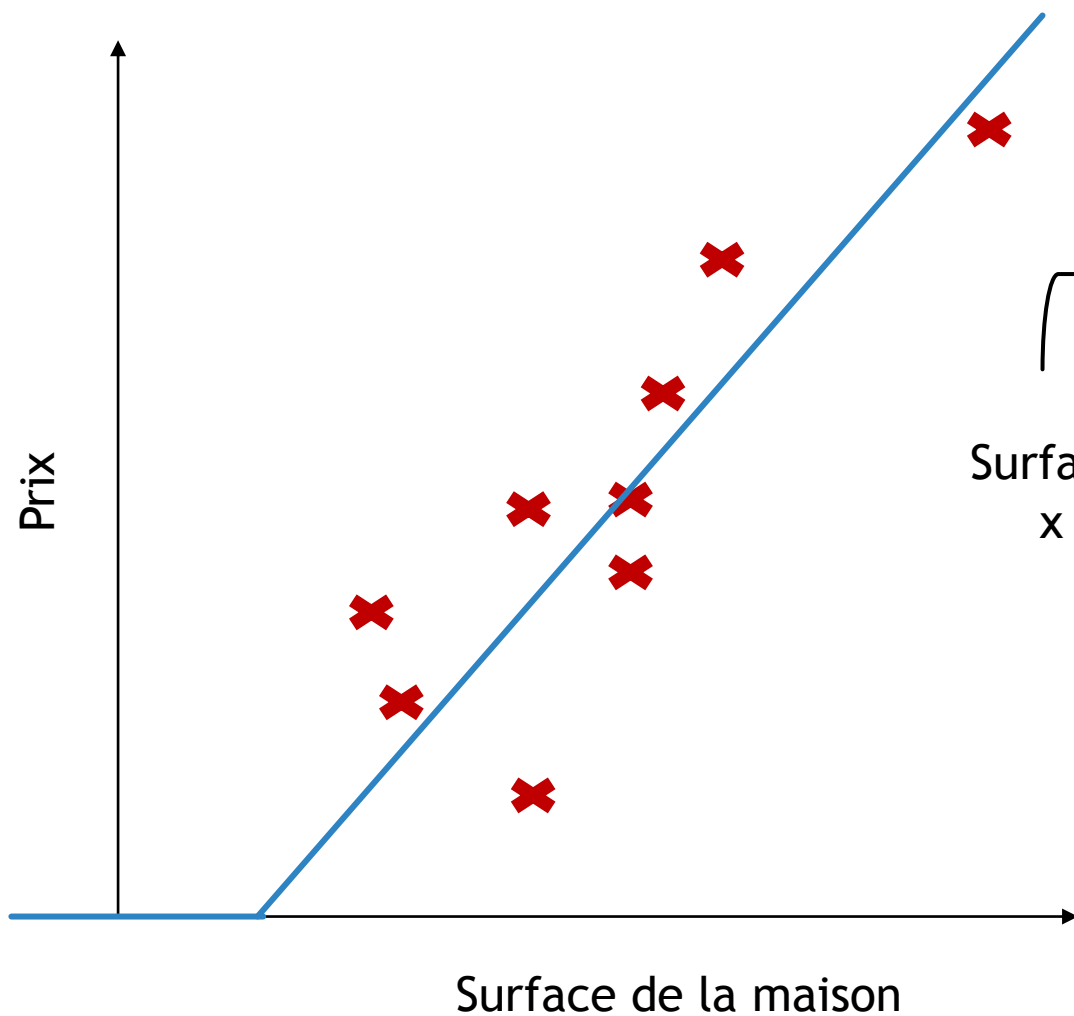
1



AI For You



# Qu'est-ce qu'un réseau neuronal ?

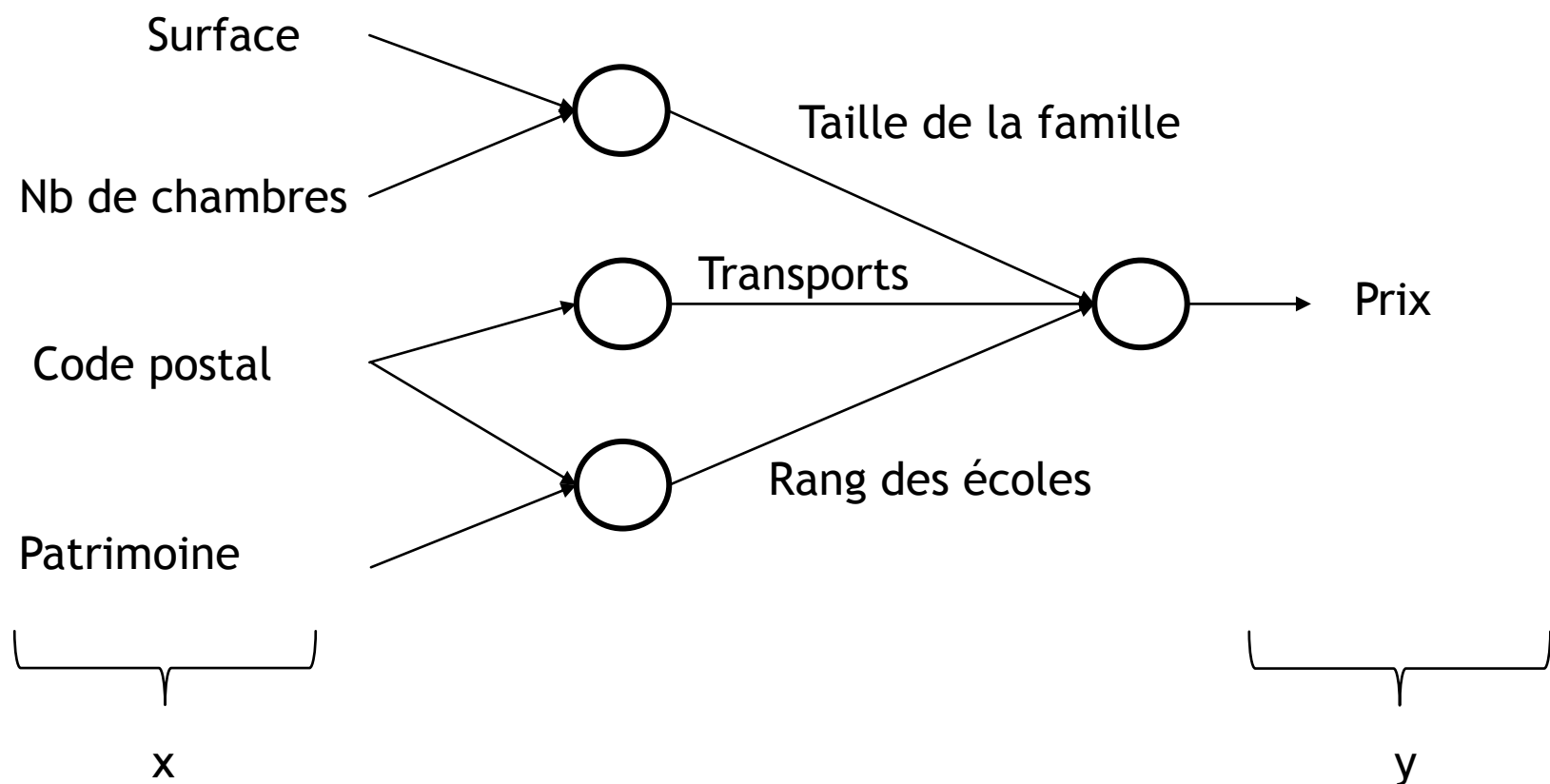


Un réseau neuronal est construit en connectant de nombreux neurones

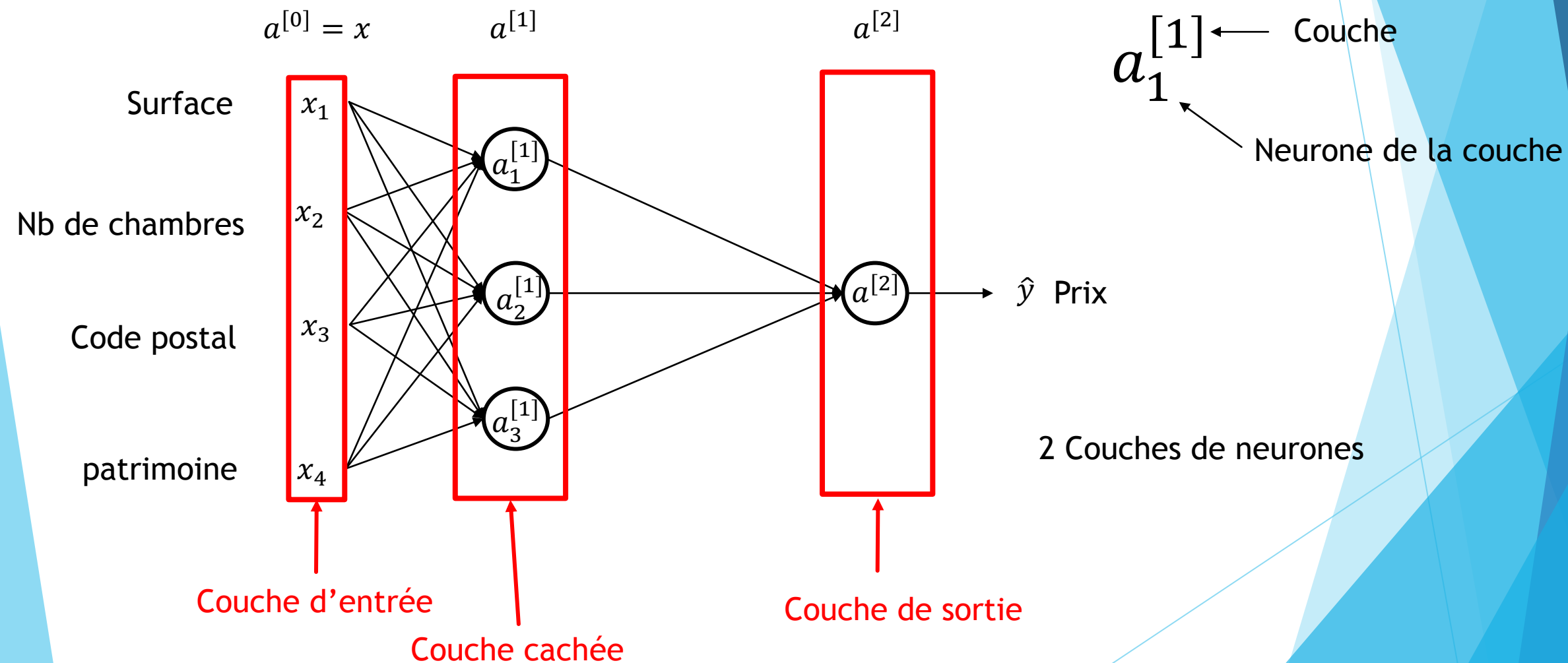
Fonction d'activation :  $g(z) = \max(0, z)$



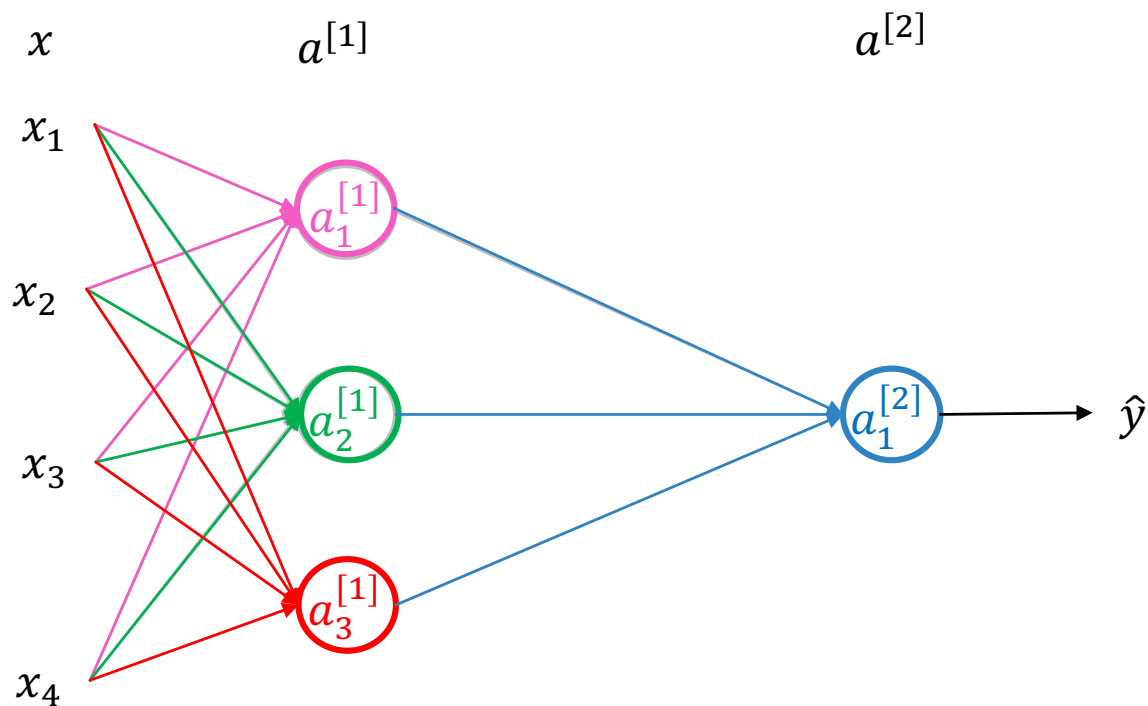
# Votre premier réseau neuronal



# Représentation des réseaux neuronaux



# Représentation des réseaux neuronaux



$$z_1^{[1]} = w_1^{[1]T} x + b_1^{[1]} ; a_1^{[1]} = g(z_1^{[1]})$$

$$z_2^{[1]} = w_2^{[1]T} x + b_2^{[1]} ; a_2^{[1]} = g(z_2^{[1]})$$

$$z_3^{[1]} = w_3^{[1]T} x + b_3^{[1]} ; a_3^{[1]} = g(z_3^{[1]})$$

$$z_1^{[2]} = w_1^{[2]T} a^{[1]} + b_1^{[2]} ; a_1^{[2]} = g(z_1^{[2]})$$



# Vectorisation

$$\begin{aligned} z_1^{[1]} &= w_1^{[1]T} x + b_1^{[1]} ; a_1^{[1]} = g(z_1^{[1]}) \\ z_2^{[1]} &= w_2^{[1]T} x + b_2^{[1]} ; a_2^{[1]} = g(z_2^{[1]}) \\ z_3^{[1]} &= w_3^{[1]T} x + b_3^{[1]} ; a_3^{[1]} = g(z_3^{[1]}) \end{aligned}$$

$$\begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \end{bmatrix} = \begin{bmatrix} w_{1,1}^{[1]} & w_{1,2}^{[1]} & w_{1,3}^{[1]} & w_{1,4}^{[1]} \\ w_{2,1}^{[1]} & w_{2,2}^{[1]} & w_{2,3}^{[1]} & w_{2,4}^{[1]} \\ w_{3,1}^{[1]} & w_{3,2}^{[1]} & w_{3,3}^{[1]} & w_{3,4}^{[1]} \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \end{bmatrix} ; \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \end{bmatrix} = g \left( \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \end{bmatrix} \right)$$

$$\begin{matrix} (3,1) & (3,4) & (4,1) & (3,1) \\ z^{[1]} & = & W^{[1]T} x & + & b^{[1]} & ; & a^{[1]} = g(z^{[1]}) \end{matrix}$$

$$\begin{matrix} (1,1) & (1,3) & (3,1) & (1,1) \\ z^{[2]} & = & w^{[2]T} a^{[1]} & + & b^{[2]} & ; & a^{[2]} = g(z^{[2]}) \end{matrix}$$



# Vectorisation sur plusieurs exemples

for ~~i = 0 to m :~~

$$a^{[1](i)} = g(w^{[1]T} x^{(i)} + b)$$

Annotations: (3,1) points to  $w^{[1]}$ , (3,4) points to  $x^{(i)}$ , (4,1) points to  $w^{[1]T}$ , (3,1) points to  $x^{(i)}$ .

$$\hat{y}^{(i)} = a^{[2](i)} = g(w^{[2]T} a^{[1](i)} + b)$$

Annotations: (1,1) points to  $w^{[2]}$ , (1,3) points to  $a^{[1](i)}$ , (3,1) points to  $w^{[2]T}$ , (1,1) points to  $a^{[1](i)}$ .

$$a^{[1]} = g(w^{[1]T} X + b)$$

Annotations: (3,m) points to  $w^{[1]}$ , (3,4) points to  $X$ , (4,m) points to  $w^{[1]T}$ , (3,1) points to  $X$ .

$$\hat{y} = a^{[2]} = g(w^{[2]T} a^{[1]} + b)$$

Annotations: (1,m) points to  $w^{[2]}$ , (1,3) points to  $a^{[1]}$ , (3,m) points to  $w^{[2]T}$ , (1,1) points to  $a^{[1]}$ .

$$x \in \mathbb{R}^{n_x} \quad x^{(i)} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix}$$

$$\hat{y} \in \mathbb{R}$$

$$X \in \mathbb{R}^{(n_x, m)} \quad X = \begin{bmatrix} x_1^{(1)} & x_1^{(2)} & x_1^{(3)} & \dots & x_1^{(m)} \\ x_2^{(1)} & x_2^{(2)} & x_2^{(3)} & \dots & x_2^{(m)} \\ x_3^{(1)} & x_3^{(2)} & x_3^{(3)} & \dots & x_3^{(m)} \\ x_4^{(1)} & x_4^{(2)} & x_4^{(3)} & \dots & x_4^{(m)} \end{bmatrix}$$

$$\hat{y} \in \mathbb{R}^m$$



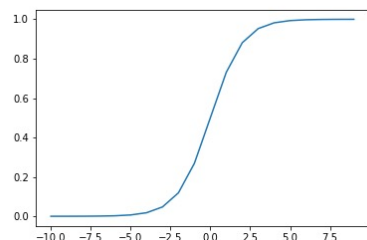
# Fonction d'activation

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

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

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

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

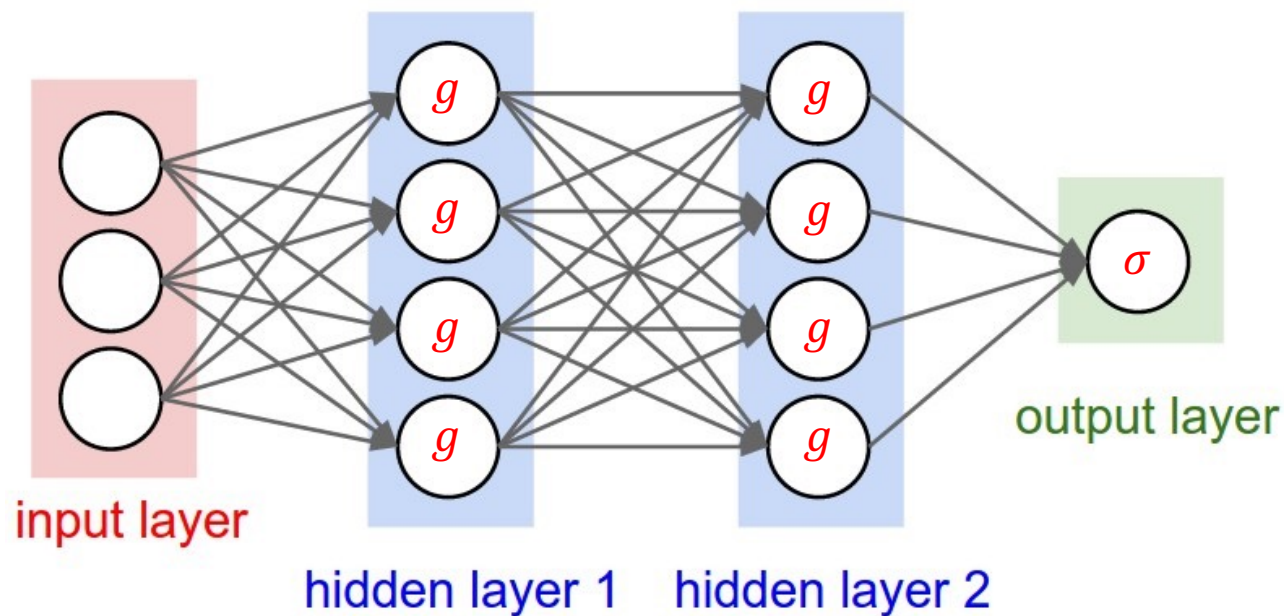


Saturation de la sigmoïde pour les grandes valeurs



Ne pas utiliser de Sigmoid dans l'activation dans un réseau de neurones

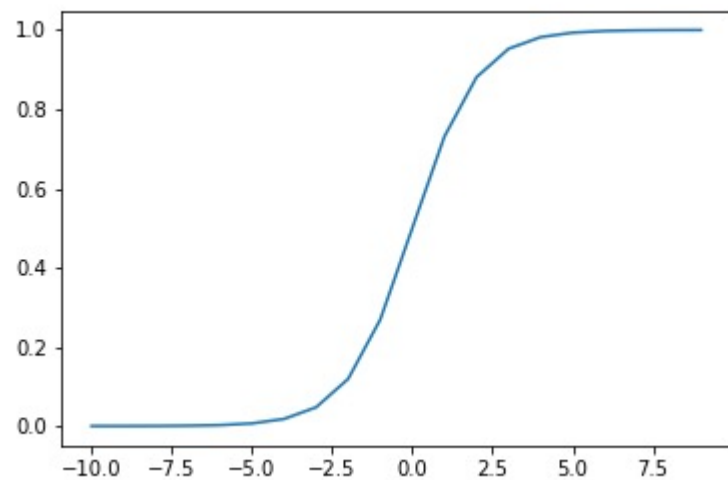
Dans le cas d'une classification binaire utilisation de la Sigmoid pour la couche de sortie







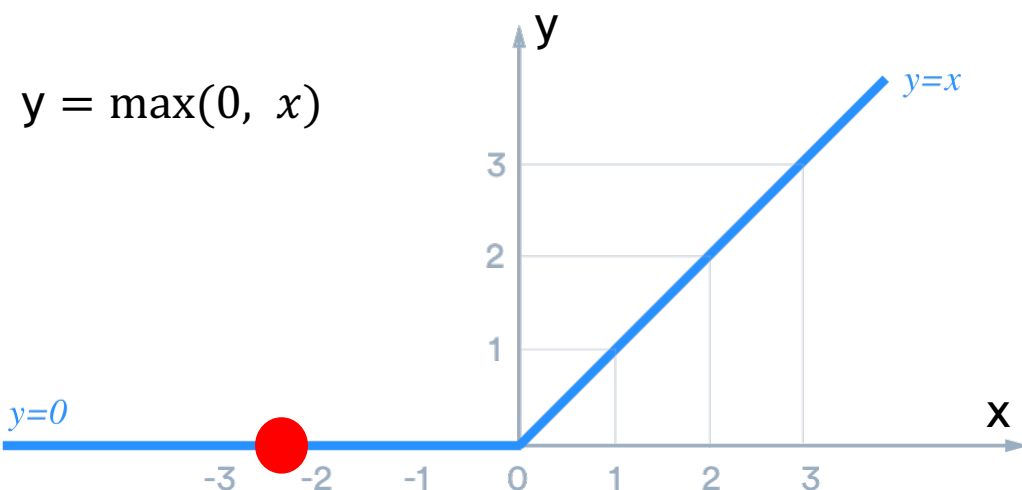
## Sigmoïd



- La sigmoïde peut saturer et conduire à des gradients évanescents.
- Non centré sur le zéro
- $e^X$  est coûteux en termes de calcul.



# Rectified Linear Unit (ReLU)



## Avantages :

ReLU ne sature pas pour les valeurs positives

ReLU est assez rapide à calculer

## Inconvénients :

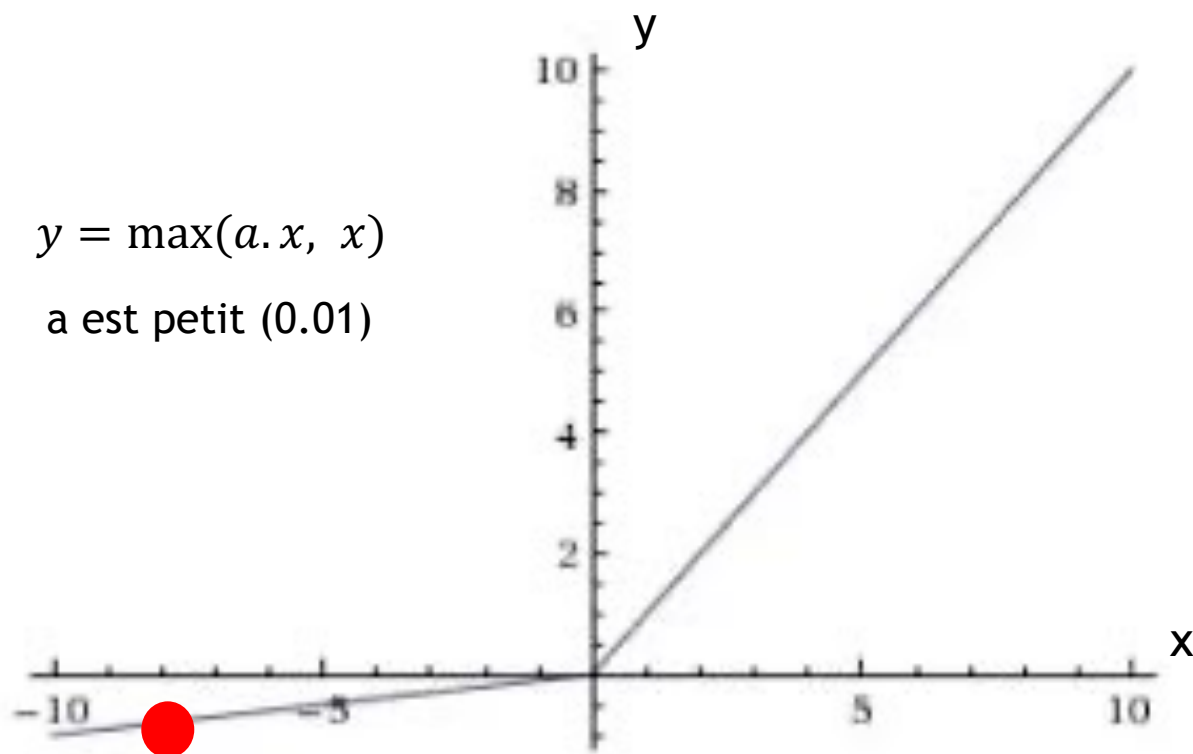
ReLU souffre d'un problème connu sous le nom de "dying ReLU".



# Leaky ReLU

$$y = \max(a \cdot x, x)$$

a est petit (0.01)

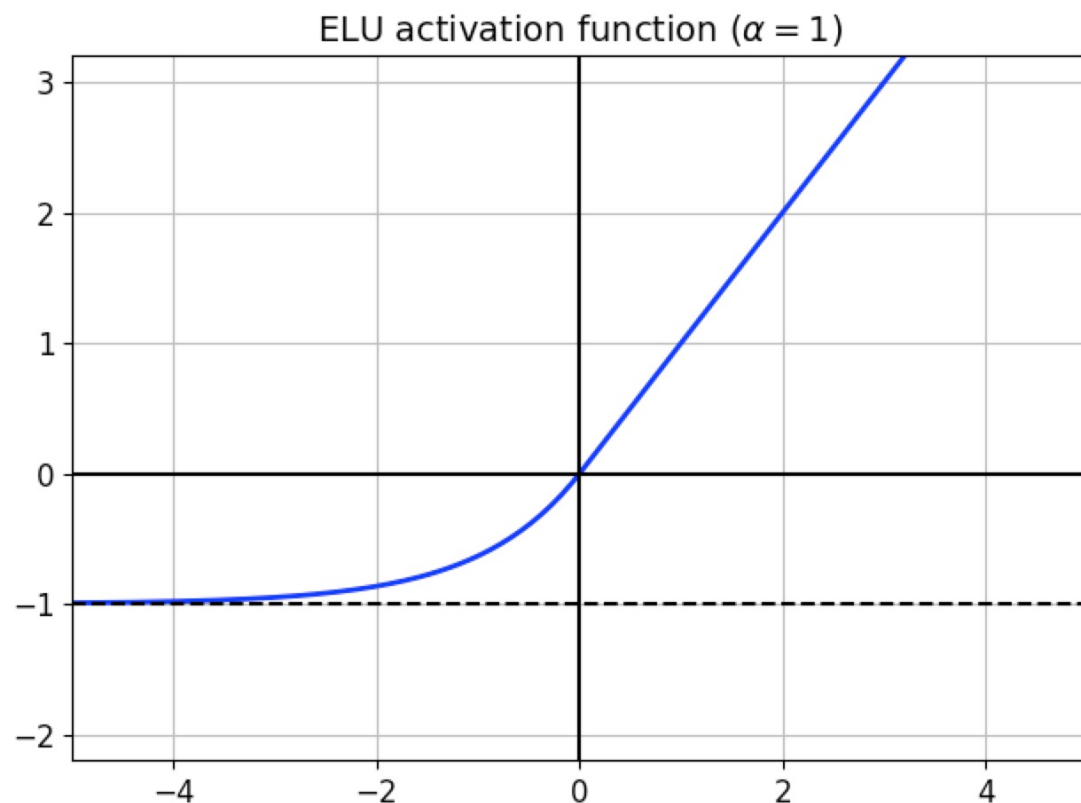


Variantes :

- Parametric leaky Relu (PReLU), si vous disposez de beaucoup de données.
- Randomized leaky Relu (RReLU), si votre réseau neuronal est surentraîné.



# Exponential Linear Unit (ELU)



$$ELU_{\alpha}(z) = \begin{cases} \alpha(\exp(z) - 1) & \text{if } z < 0 \\ z & \text{if } z \geq 0 \end{cases}$$

## Avantages :

- Il prend des valeurs négatives, ce qui permet au neurone d'avoir une moyenne plus proche de 0
- Il a un gradient non nul pour  $z < 0$ , ce qui évite le problème des 'dying ReLU'.

## Inconvénients

- Il est plus lent à calculer que ReLU.



# Fonctions d'activations conseils

- ▶ En général pour la performance :

ELU > Leaky ReLU > ReLU

- ▶ Si vous vous souciez du temps d'exécution :

Leaky ReLU > ELU

- ▶ Par défaut :

ReLU



## Descente de gradient pour les réseaux de neurones

Paramètres :  $w^{[1]}$ ,  $b^{[1]}$ ,  $w^{[2]}$ ,  $b^{[2]}$        $n_x = n^{[0]}$ ,  $n^{[1]}$ ,  $n^{[2]} = 1$

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

*Note: An arrow points from  $a^{[2]}$  to  $\hat{y}$  in the loss function.*

Descente de gradient :

Répéter {

Calculer la prédiction ( $\hat{y}^{(i)}, i = (1, \dots, m)$ )

$$\boxed{dw^{[1]}} = \frac{dJ}{dw^{[1]}}, \boxed{db^{[1]}} = \frac{dJ}{db^{[1]}}, \boxed{dw^{[2]}} = \frac{dJ}{dw^{[2]}}, \boxed{db^{[2]}} = \frac{dJ}{db^{[2]}}$$

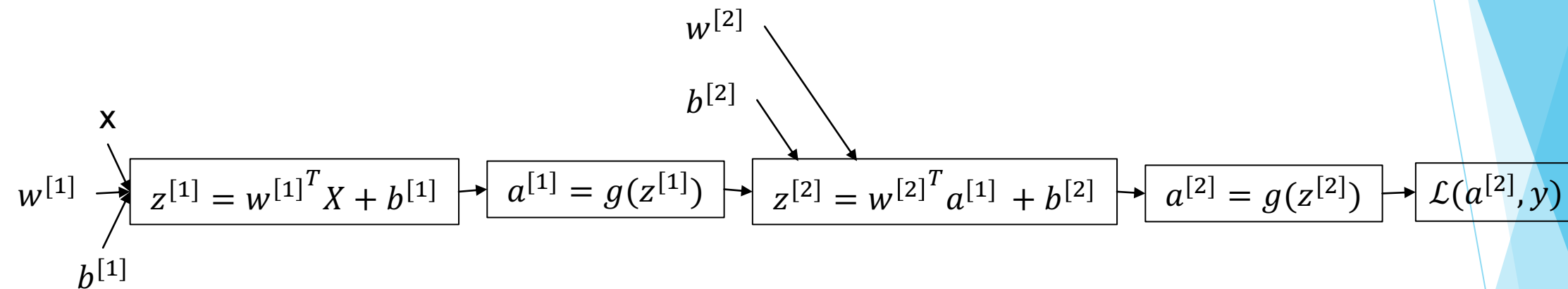
$$w^{[1]} := w^{[1]} - \alpha dw^{[1]}$$

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

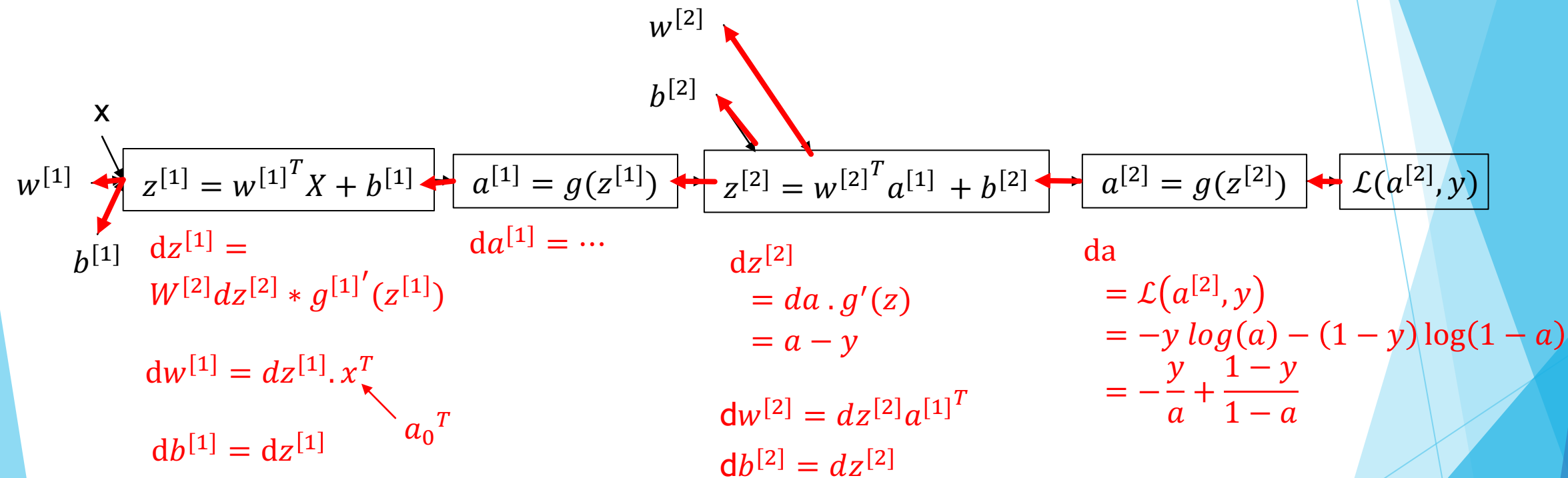
$$w^{[2]} := w^{[2]} - \alpha dw^{[2]}$$

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

# Forward propagation



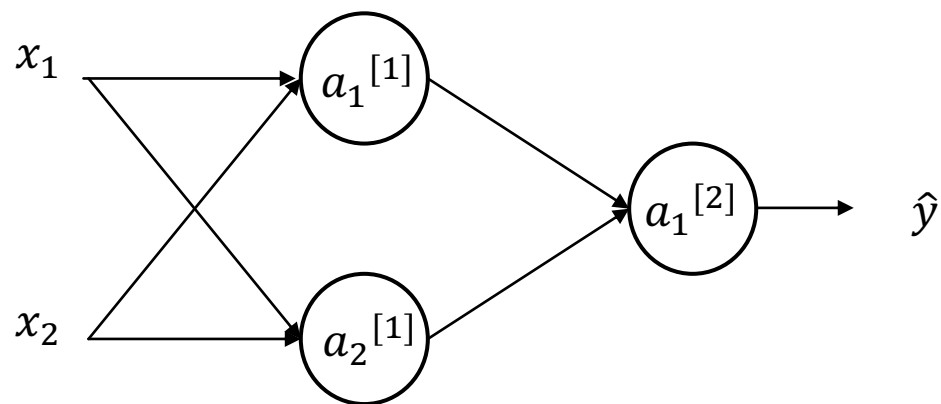
# Back propagation







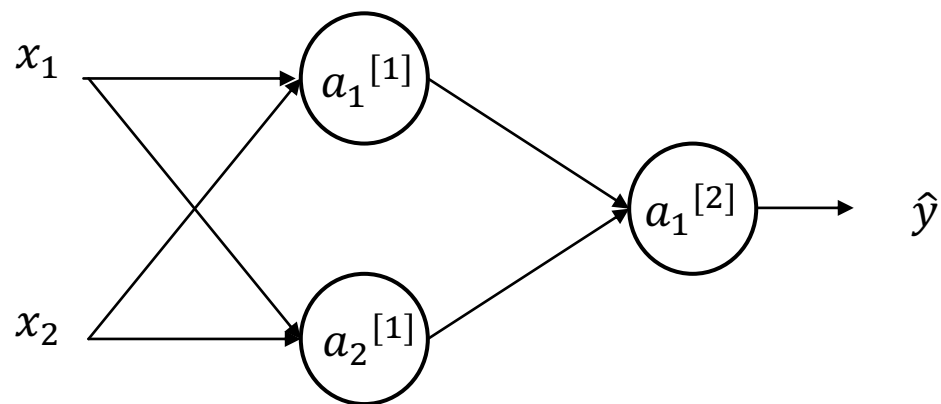
# Initialiser les poids à zéro ?



$$w^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \quad b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$w^{[1]} = \begin{bmatrix} p & l \\ p & l \end{bmatrix} \leftarrow a_1^{[1]} = a_2^{[1]} \longrightarrow dz_1^{[1]} = dz_2^{[1]} \longrightarrow dw^{[1]} = \begin{bmatrix} u & v \\ u & v \end{bmatrix} \longrightarrow w^{[1]} := w^{[1]} - \alpha dw$$

# Initialiser le poids de façon aléatoire



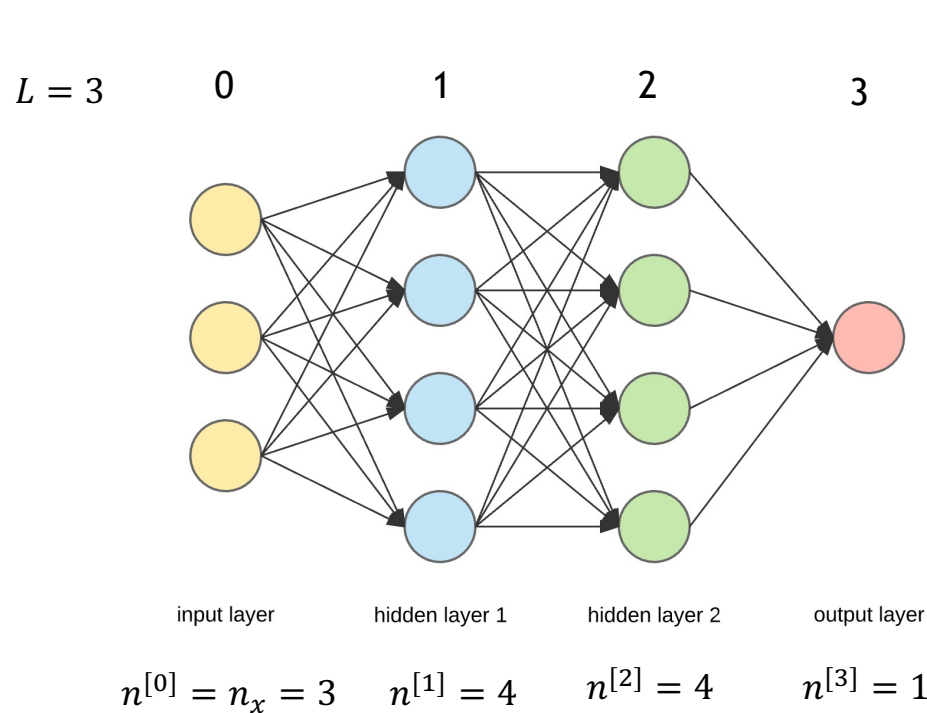
$$w^{[1]} = np.random.rand((2,2)) \cdot \text{Petit nombre}$$

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

$$w^{[1]} = np.random.rand((1,2)) \cdot \text{Petit nombre}$$

$$b^{[1]} = 0$$

# Bien choisir les dimensions de la matrice



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

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

$$Z^{[1]}_{(n^{[1]}, m)} = w^{[1]T}_{(n^{[1]}, n^{[0]})} A^{[0]}_{(n^{[0]}, m)} + b^{[1]}_{(n^{[1]}, 1)}$$

$\swarrow$   
 $x$

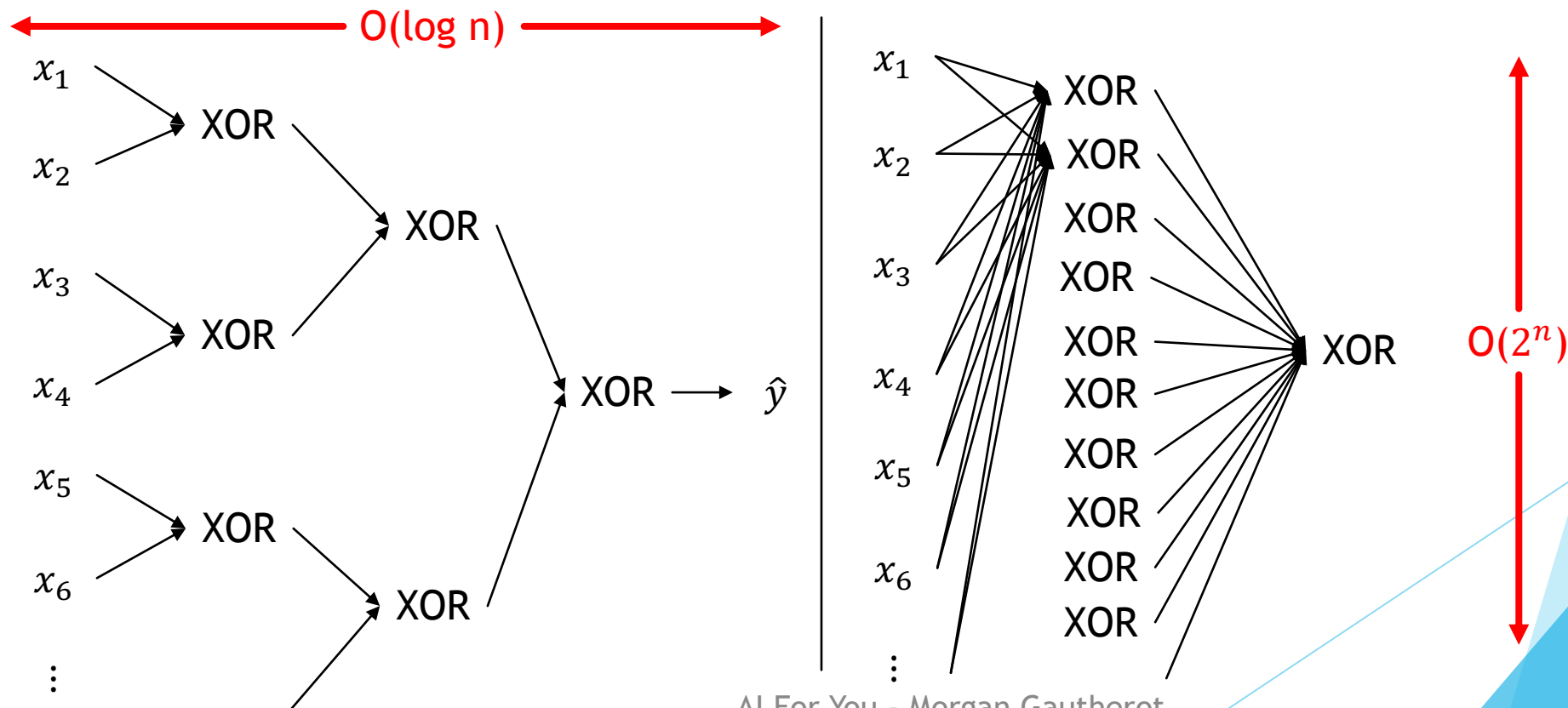
$$Z^{[2]}_{(n^{[2]}, m)} = w^{[2]T}_{(n^{[2]}, n^{[1]})} A^{[1]}_{(n^{[1]}, m)} + b^{[2]}_{(n^{[2]}, 1)}$$

$$Z^{[3]}_{(n^{[3]}, m)} = w^{[3]T}_{(n^{[3]}, n^{[2]})} A^{[2]}_{(n^{[2]}, m)} + b^{[3]}_{(n^{[3]}, 1)}$$

# Pourquoi des représentations profondes ?

Pour approximer une fonction, une architecture de réseau neuronal plus profonde nécessite moins de paramètres pour être entraînée qu'une architecture moins profonde.

$$y = x_1 \text{ XOR } x_2 \text{ XOR } x_3 \text{ XOR } x_4 \text{ XOR } x_5 \text{ XOR } x_6 \text{ XOR } \dots$$



# Résumé de l'apprentissage profond

