

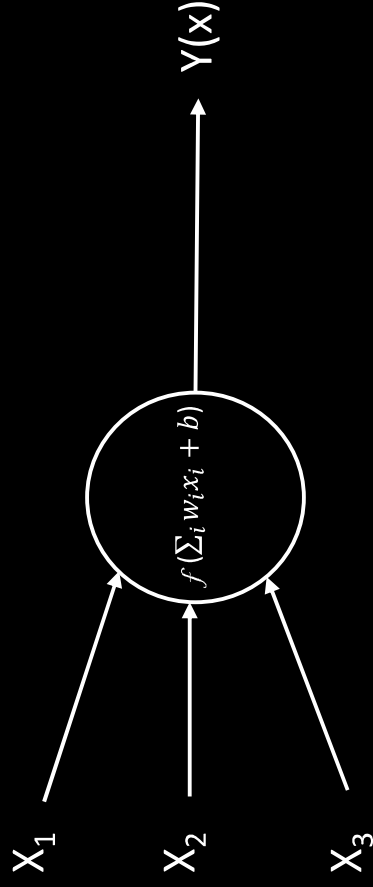
Evolution du Deep Learning



Le Perceptron

1. Neurone Formel et le perceptron Multicouche.
2. Fonction d'activation, en particulier la fonction sigmoïde (Logistic)
3. Fonction Cout et minimisation de l'erreur.
4. La Descente du gradient.
5. Algorithme de la descente de gradient.

Neurone formel

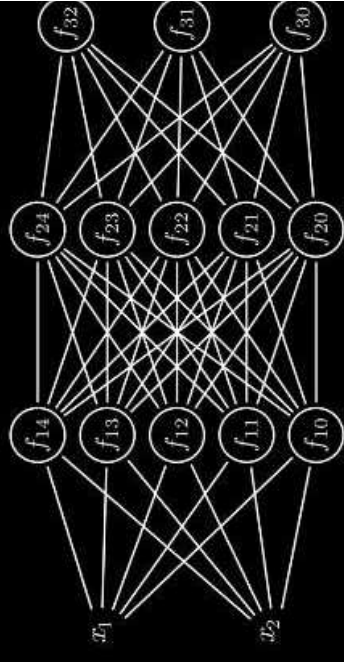


x_i : Entrées

f : Fonction d'activation

Y : Sortie

Perceptron multicouche



Entrées

f_{ij} : Fonction d'activation

Fonction d'activation

Sigmoid

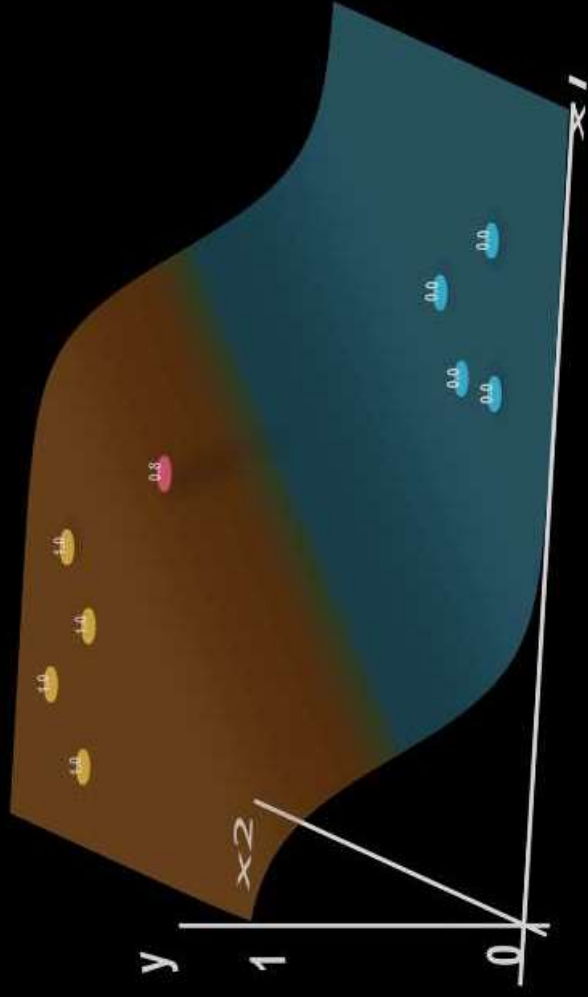
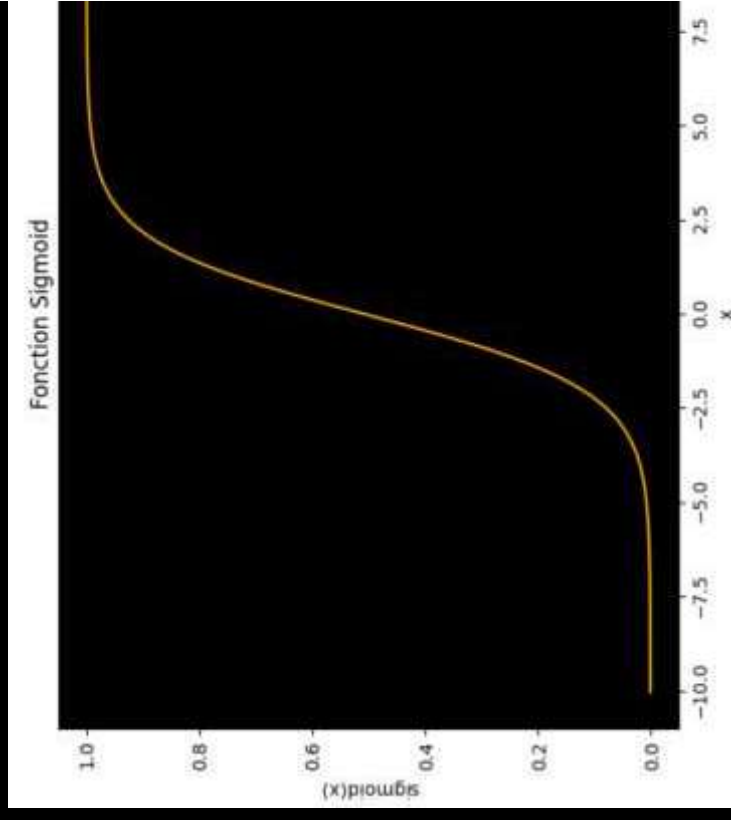
$$f(x) = \frac{1}{1 + e^{-x}}$$

Relu

$$f(x) = \max(0, x)$$

Heaviside

$$H(x) = \begin{cases} 1, & \text{Si } x > 0 \\ \frac{1}{2}, & \text{Si } x = 0 \\ 0, & \text{Si } x < 0 \end{cases}$$



Source : Machine Learnia

Minimisation de l'erreur

Pour minimiser les erreurs de notre modèle, nous utiliserons la fonction de coût **Log Loss**:

$$-\frac{1}{m} * \sum_{i=1}^m y_i * \log(a_i) + (1 - y_i) * \log(1 - a_i)$$

$\begin{cases} a_i: & \text{la probabilité d'être dans la classe } i \\ y_i: & \text{la sortie } i \text{ de l'échantillon} \end{cases}$

Cette formule est dérivée de la définition statistique de **Log-vraisemblance** :

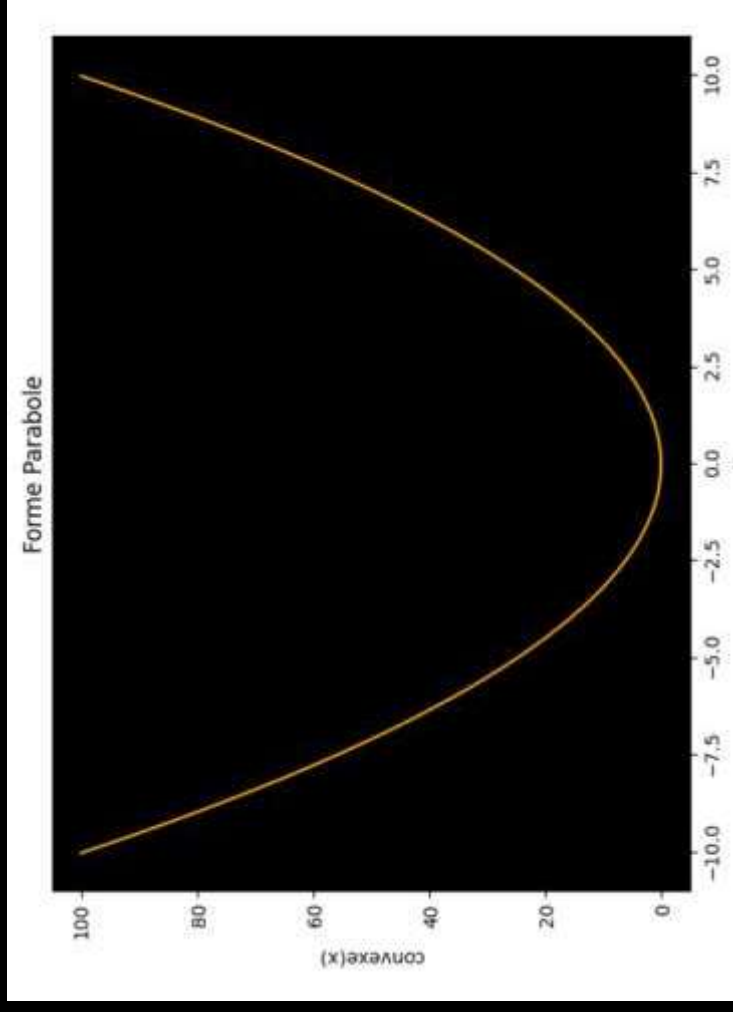
$$\log\left(\prod_{i=1}^m P(X = x_i)\right) = \log\left(\prod_{i=1}^m a_i^{x_i} * (1 - a_i)^{1-x_i}\right)$$

avec $x_i \in \{0,1\}$

En effet, pour maximiser la vraisemblance, nous minimisons la fonction de coût car il existe de nombreux algorithmes d'optimisation pour minimiser les fonctions.

L' algorithme de la descente du gradient

Il s'agit d'un algorithme de minimisation qui consiste à ajuster les paramètres **W** et le biais **b** afin de minimiser le coût. On calcule ainsi le gradient de la fonction coût



- Si $\frac{\partial F}{\partial x} > 0$, la fonction F augmente lorsque x augmente
 - Sinon, la fonction F diminue lorsque x augmente
- Dans notre cas, nous chercherons à atteindre un minimum global pour la fonction de coût. Nous déterminons donc la fonction par rapport aux paramètres, et nous augmenterons la valeur de ces paramètres en fonction du signe de la dérivée jusqu'à atteindre le minimum.

La formule conçue pour ce travail est :

$$W_{i+1} = W_i - \mu * \frac{\partial L}{\partial W_i}$$

En résumé

Neurone Formel

$$z = \sum_{i=1}^m w_i x_i + w_0$$

$$y(\text{sortie}) = \begin{cases} 0, & \text{Si } z < 0 \\ 1, & \text{Si } z \geq 0 \end{cases}$$

Perceptron

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

$$-\frac{1}{m} * \sum_{i=1}^m y_i * \log(a_i) + (1 - y_i) * \log(1 - a_i)$$

$$W_{i+1} = W_i - \mu * \frac{\partial L}{\partial W_i}$$

Vectorisation des équations pour des données larges

1) Vectorisation du modèle:

$$\left\{ \begin{array}{l} X = \begin{bmatrix} x_1^{(1)} & \dots & x_n^{(1)} \\ \vdots & & \vdots \\ x_1^{(m)} & \dots & x_n^{(m)} \end{bmatrix} \\ W = \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix} \\ B = \begin{bmatrix} b \\ \vdots \\ b \end{bmatrix} \end{array} \right\} \quad \Rightarrow \quad Z = \begin{bmatrix} z^{(1)} \\ \vdots \\ z^{(m)} \end{bmatrix} = \begin{bmatrix} w_1 * x_1^{(1)} + w_2 * x_2^{(1)} + \dots + w_n * x_n^{(1)} + b \\ w_1 * x_1^{(m)} + w_2 * x_2^{(m)} + \dots + w_n * x_n^{(m)} + b \end{bmatrix} =$$

2) Vectorisation de la fonction d'activation:

$$A = \begin{bmatrix} \sigma(z^{(1)}) \\ \vdots \\ \sigma(z^{(m)}) \end{bmatrix} = \sigma \left(\begin{bmatrix} z^{(1)} \\ \vdots \\ z^{(m)} \end{bmatrix} \right) = \sigma(Z)$$

$$\text{avec : } \sigma(z) = \frac{1}{1 + e^{-z}} \quad (\text{fonction log})$$

3) Vectorisation de la fonction cout:

Rappelons que :

$$L = -\frac{1}{m} * \sum_{i=1}^m y_i * \log(a_i) + (1 - y_i) * \log(1 - a_i)$$

$$L = -\frac{1}{m} * \text{Somme_des_lignes}(y * \log(A) + (1 - y) * \log(1 - A))$$

*Produit terme à terme

$$\left\{ \begin{array}{l} y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix} \\ A = \begin{bmatrix} \sigma(z^{(1)}) \\ \vdots \\ \sigma(z^{(m)}) \end{bmatrix} \end{array} \right\}, \text{ avec}$$

4) Vectorisation de l'algorithme de la descente de gradient:

$$\left\{ \begin{array}{l} W = \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix} \\ B = \begin{bmatrix} b \\ \vdots \\ b \end{bmatrix} \\ \frac{\partial L}{\partial W} = \begin{bmatrix} \frac{\partial L}{\partial w_1} \\ \vdots \\ \frac{\partial L}{\partial w_n} \end{bmatrix} \\ \frac{\partial L}{\partial B} = \begin{bmatrix} \frac{\partial L}{\partial b} \\ \vdots \\ \frac{\partial L}{\partial b} \end{bmatrix} \end{array} \right\}$$

$$W = W - \mu * \frac{\partial L}{\partial W}$$

$$B = B - \mu * \frac{\partial L}{\partial B}$$

Calcul des gradients

1) Calcul de $\frac{\partial L}{\partial w_i}$:

$$\text{On a } L = -\frac{1}{m} * \sum_{i=1}^m y_i * \log(a_i) + (1 - y_i) * \log(1 - a_i)$$

$$\text{alors } L = -\frac{1}{m} * \sum_{i=1}^m y_i * \log\left(\frac{1}{1+e^{-z_i}}\right) + (1 - y_i) * \log\left(1 - \frac{1}{1+e^{-z_i}}\right)$$

$$\text{alors } L = -\frac{1}{m} * \sum_{i=1}^m (-y_i * \log(1 + e^{-z_i})) + (1 - y_i) * (-\log(1 + e^{-z_i}) - z_i)$$

$$\text{Donc } L = -\frac{1}{m} * \sum_{i=1}^m -\log(e^{z_i} + 1) + y_i z_i$$

$$\frac{\partial L}{\partial w_j} = -\frac{1}{m} * \sum_{i=1}^m (y_i - a_i) \times x_j^{(i)}$$

Finalement :

2) Calcul de $\frac{\partial L}{\partial b}$:

De façon analogue :

$$\frac{\partial L}{\partial b} = -\frac{1}{m} * \sum_{i=1}^m (y_i - a_i)$$

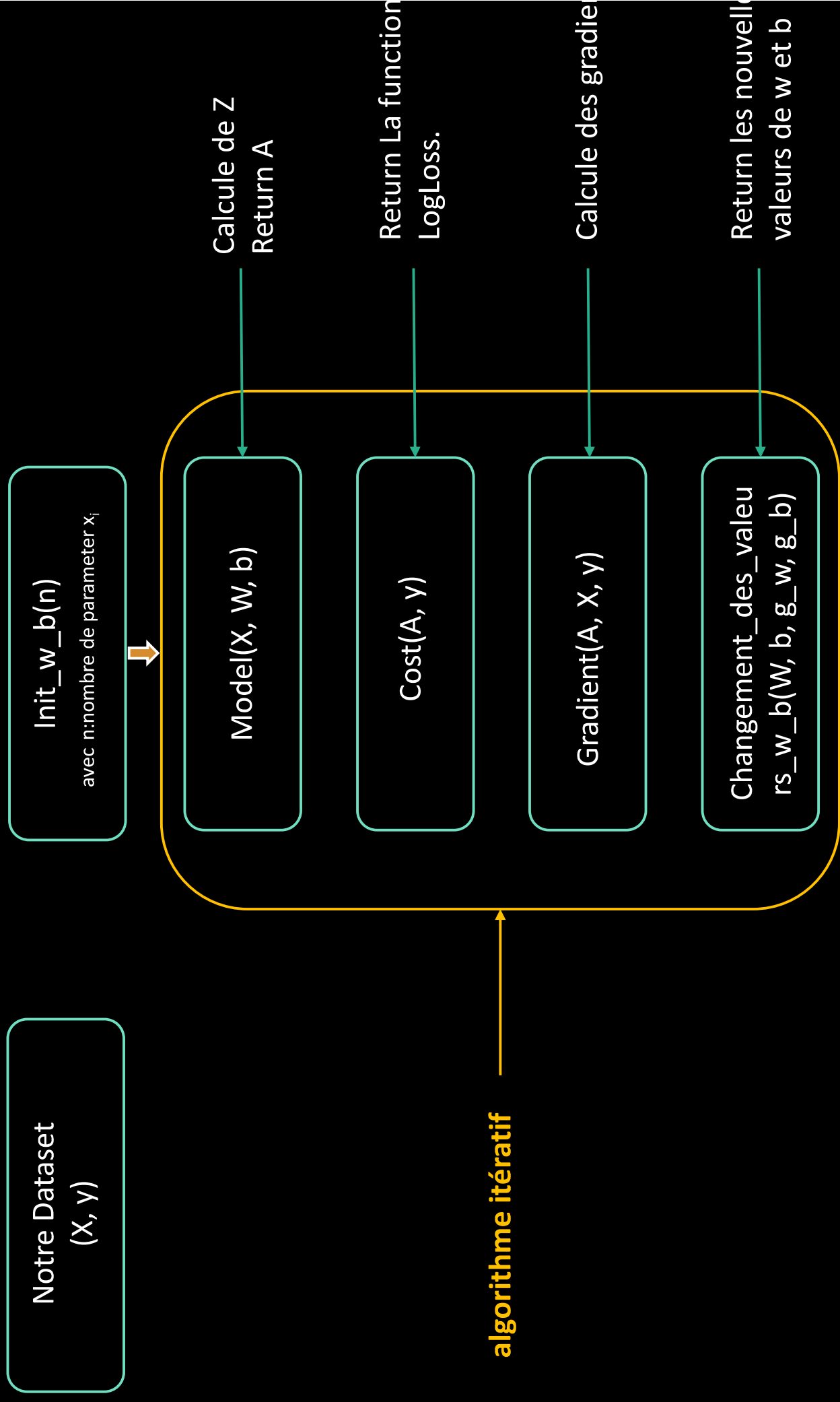
Ainsi sous une forme vectorisée :

$$\frac{\partial L}{\partial W} = \begin{bmatrix} \frac{\partial L}{\partial w_1} \\ \vdots \\ \frac{\partial L}{\partial w_n} \end{bmatrix} = -\frac{1}{m} \begin{bmatrix} \sum_{i=1}^m (y_i - a_i) \times x_1^{(i)} \\ \vdots \\ \sum_{i=1}^m (y_i - a_i) \times x_n^{(i)} \end{bmatrix} = -\frac{1}{m} \begin{bmatrix} x_1^{(1)} \\ \vdots \\ x_n^{(1)} \end{bmatrix} \cdot \begin{pmatrix} x_1^{(m)} \\ \vdots \\ x_n^{(m)} \end{pmatrix} \cdot \left(\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} - \begin{bmatrix} a_1 \\ \vdots \\ a_n \end{bmatrix} \right) = -\frac{1}{m} X^T \cdot (y - A)$$

De meme:

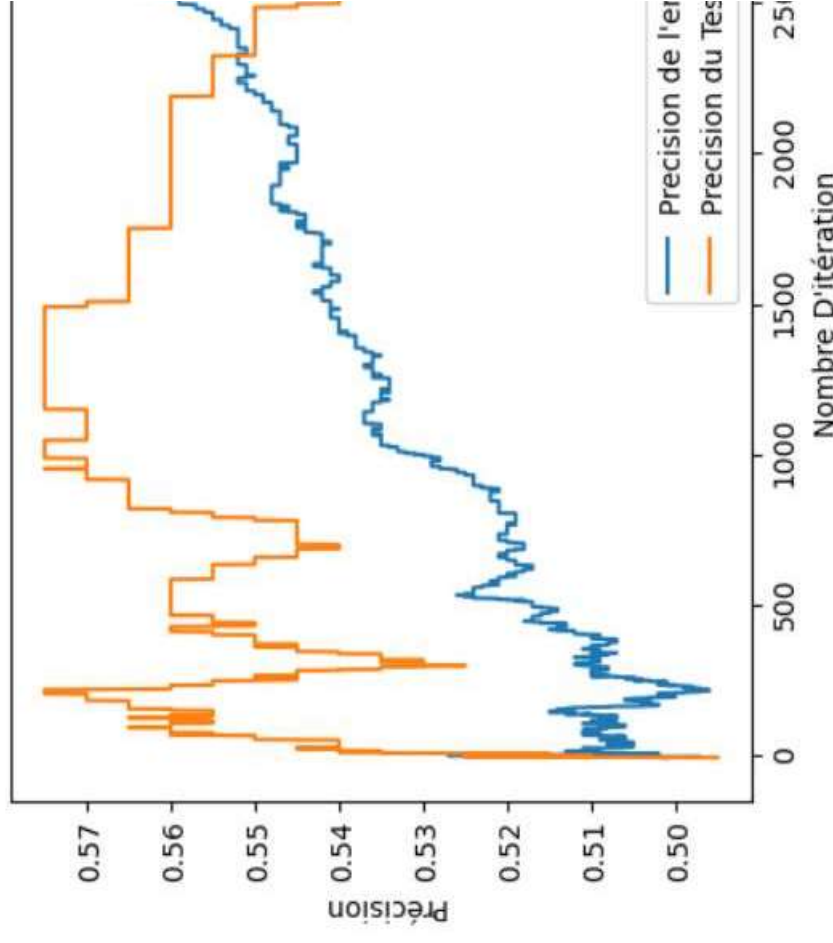
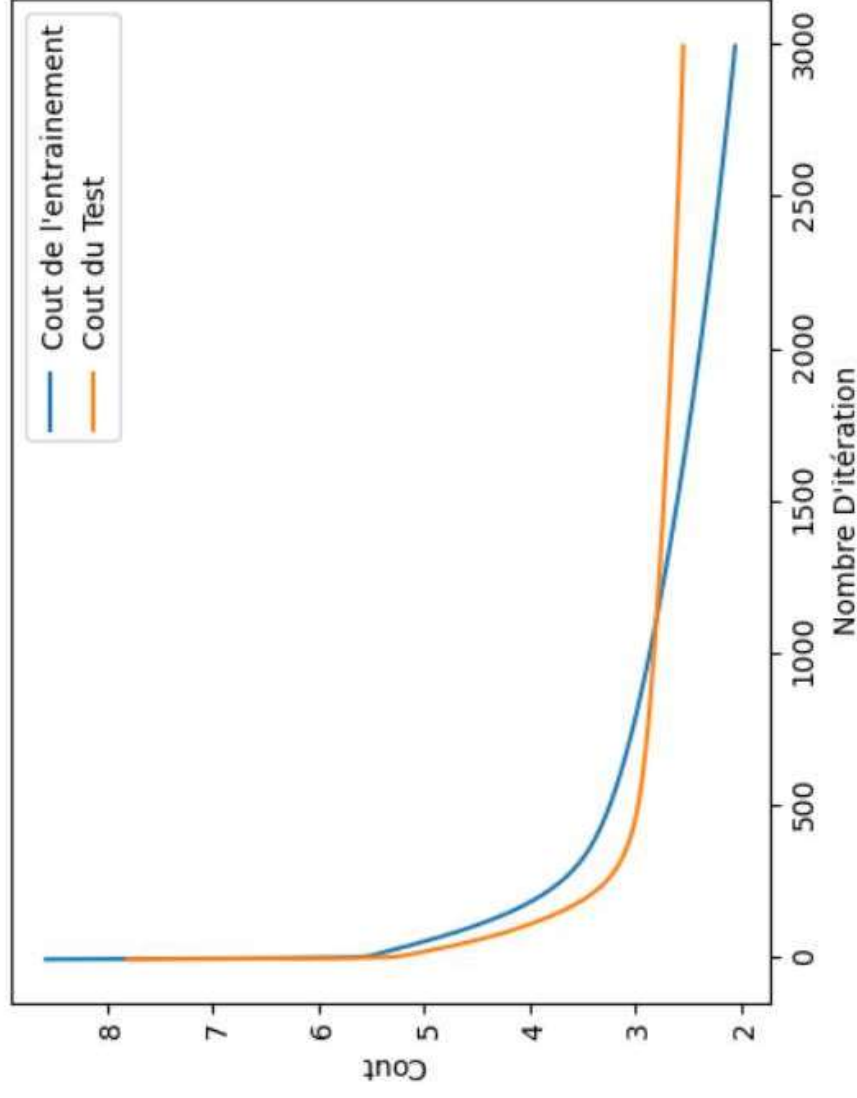
$$\frac{\partial L}{\partial B} = \begin{bmatrix} \frac{\partial L}{\partial b} \\ \vdots \\ \frac{\partial L}{\partial b} \end{bmatrix} = -\frac{1}{m} \begin{bmatrix} \sum_{i=1}^m (y_i - a_i) \\ \vdots \\ \sum_{i=1}^m (y_i - a_i) \end{bmatrix} = -\frac{1}{m} \begin{bmatrix} \sum (y - A) \\ \vdots \\ \sum (y - A) \end{bmatrix}$$

Programmation du neurone artificiel



Application : Chat ou Chien ?

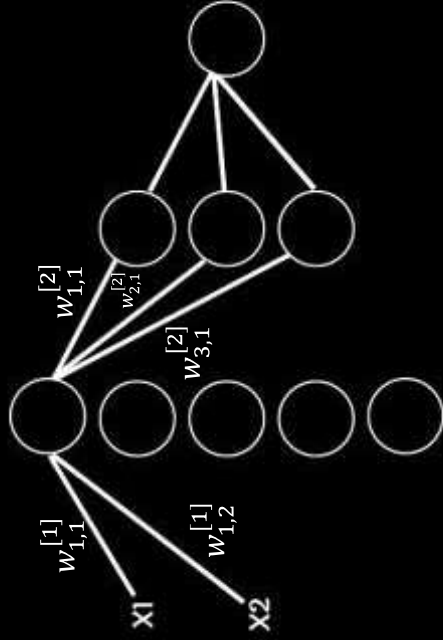
Résultat:



Vers les réseaux de neurones

L'insuffisance du perceptron tout seul nous oblige à travailler avec plusieurs neurones connectés entre eux pour faire des tâches plus complexes et plus précises :

1) Forward Propagation:



Notation:

$W_{i,j}^{[\alpha]}$: poids associé au neurone i
et provenant de l'entrée j
et α le numero de la couche

Véctorisation des couches d'un réseau de neurone:

Fonction d'agrégation:

$$\begin{cases} Z^{[1]} = W^{[1]} \cdot X + B^{[1]} \\ Z^{[\alpha]} = W^{[\alpha]} \cdot A^{[\alpha-1]} + B^{[\alpha]} \end{cases}$$

Fonction d'activation:

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

2) Back Propagation:

On procède de la même façon que pour un perceptron, sauf qu'ici pour actualiser le poids de connexion $w_{j,q}^h$ du neurone j de la couche $h-1$ vers le neurone q de la couche h , nous devons calculer $\frac{\partial L}{\partial w_{j,q}^h}$. Pour ce faire, nous allons appliquer le théorème de dérivation des fonctions composées (Règle de la chaîne).

D'après les calculs faits sur la fonction coût d'un perceptron, on obtient :

Cas de $C = 2$:

$$\begin{aligned}\frac{\partial L}{\partial w_{j,q}^{[2]}} &= \frac{\partial L}{\partial A^{[2]}} \times \frac{\partial A^{[2]}}{\partial Z^{[2]}} \times \frac{\partial Z^{[2]}}{\partial w_{j,q}^{[2]}} \\ \frac{\partial L}{\partial B^{[2]}} &= \frac{\partial L}{\partial A^{[2]}} \times \frac{\partial A^{[2]}}{\partial Z^{[2]}} \times \frac{\partial Z^{[2]}}{\partial B^{[2]}} \\ \frac{\partial L}{\partial w_{j,q}^{[1]}} &= \frac{\partial L}{\partial A^{[2]}} \times \frac{\partial A^{[2]}}{\partial Z^{[2]}} \times \frac{\partial Z^{[2]}}{\partial A^{[1]}} \times \frac{\partial A^{[1]}}{\partial Z^{[1]}} \times \frac{\partial Z^{[1]}}{\partial w_{j,q}^{[1]}} \\ \frac{\partial L}{\partial B^{[1]}} &= \frac{\partial L}{\partial A^{[2]}} \times \frac{\partial A^{[2]}}{\partial Z^{[2]}} \times \frac{\partial Z^{[2]}}{\partial A^{[1]}} \times \frac{\partial A^{[1]}}{\partial Z^{[1]}} \times \frac{\partial Z^{[1]}}{\partial B^{[1]}}\end{aligned}$$

Avec:

$$\begin{aligned}\frac{\partial L}{\partial A} &= \sum \frac{-y}{A} + \frac{1-y}{A} \\ \frac{\partial A}{\partial Z} &= A \times (1-A) \\ \frac{\partial Z}{\partial W} &= A\end{aligned}$$

On constate alors une répétition de calculs, on utilisera alors la mémoïsation pour optimiser nos calculs.

Notation:

$[C]$: numéro de la couche.

Notation δ pour les valeurs de la fonction de transfert.

On pose :

$$base^{[2]} = \frac{\partial L}{\partial A^{[2]}} \times \frac{\partial A^{[2]}}{\partial Z^{[2]}}$$

$$base^{[1]} = \frac{\partial L}{\partial A^{[2]}} \times \frac{\partial A^{[2]}}{\partial Z^{[2]}} \times \frac{\partial Z^{[2]}}{\partial A^{[1]}} \times \frac{\partial A^{[1]}}{\partial Z^{[1]}} = base^{[2]} \times \frac{\partial Z^{[2]}}{\partial A^{[1]}} \times \frac{\partial A^{[1]}}{\partial Z^{[1]}} = W^{[2]^T} \times base^{[2]} \times A^{[1]} \times (1 -$$

Cas général:

En prenant en considération la technique de mémorisation, on obtient les équations suivantes pour une c

$$\frac{\partial L}{\partial W^{[c]}} = \frac{1}{m} \times base^{[c]} . A^{[c-1]^T} : \text{la somme est incluse dans le produit matricielle}$$

$$\frac{\partial L}{\partial B^{[c]}} = \frac{1}{m} \times base^{[c]}$$

$$base^{[c-1]} = W^{[c]^T} \times base^{[c]} \times A^{[c-1]} \times (1 - A^{[c-1]}) : \text{technique de mémorisation.}$$

Convolution

1)Operation de convolution: (Padding = 'same' ou 'valid', Stride = (x, y), kernel = taille du filtre)

0	0	0	0	0
0	2	50	20	0
0	20	18	20	0
0	20	20	20	0

*

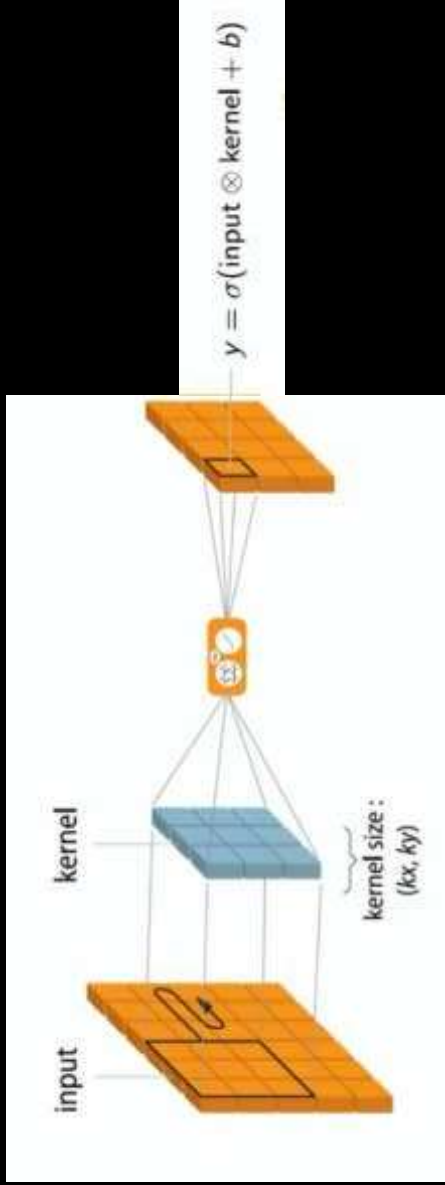
-1	2	2
2	6	2
2	2	-1

=

134			

Portion d'une image

Filtre (Kernel)



2)Neurone de Convolution:

Maxpooling

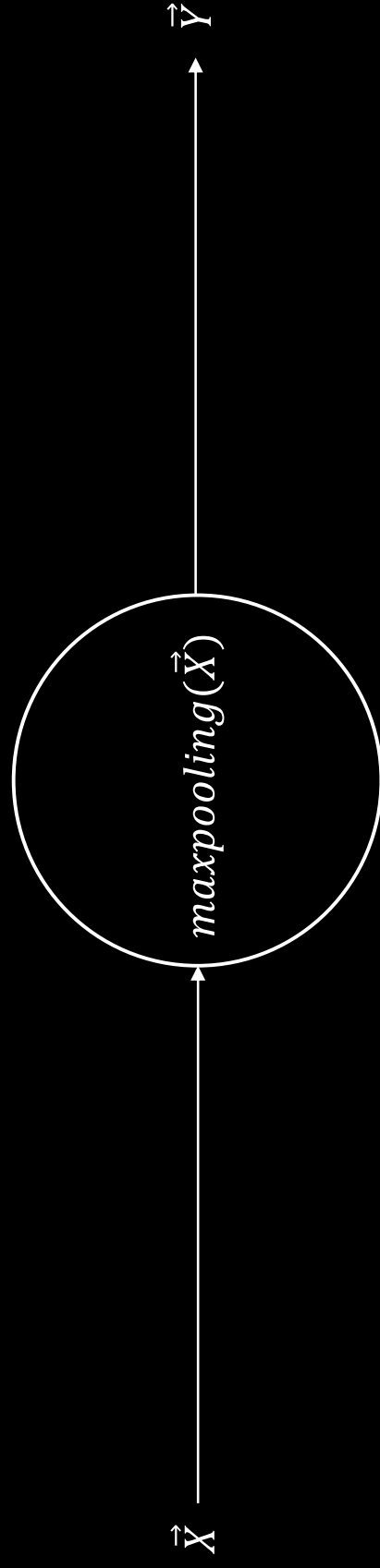
1) Operation de maxpooling: (Stride = (x, y), kernel = taille du filtre)

0	0	0	0
0	2	50	20
0	20	18	20
0	20	20	20

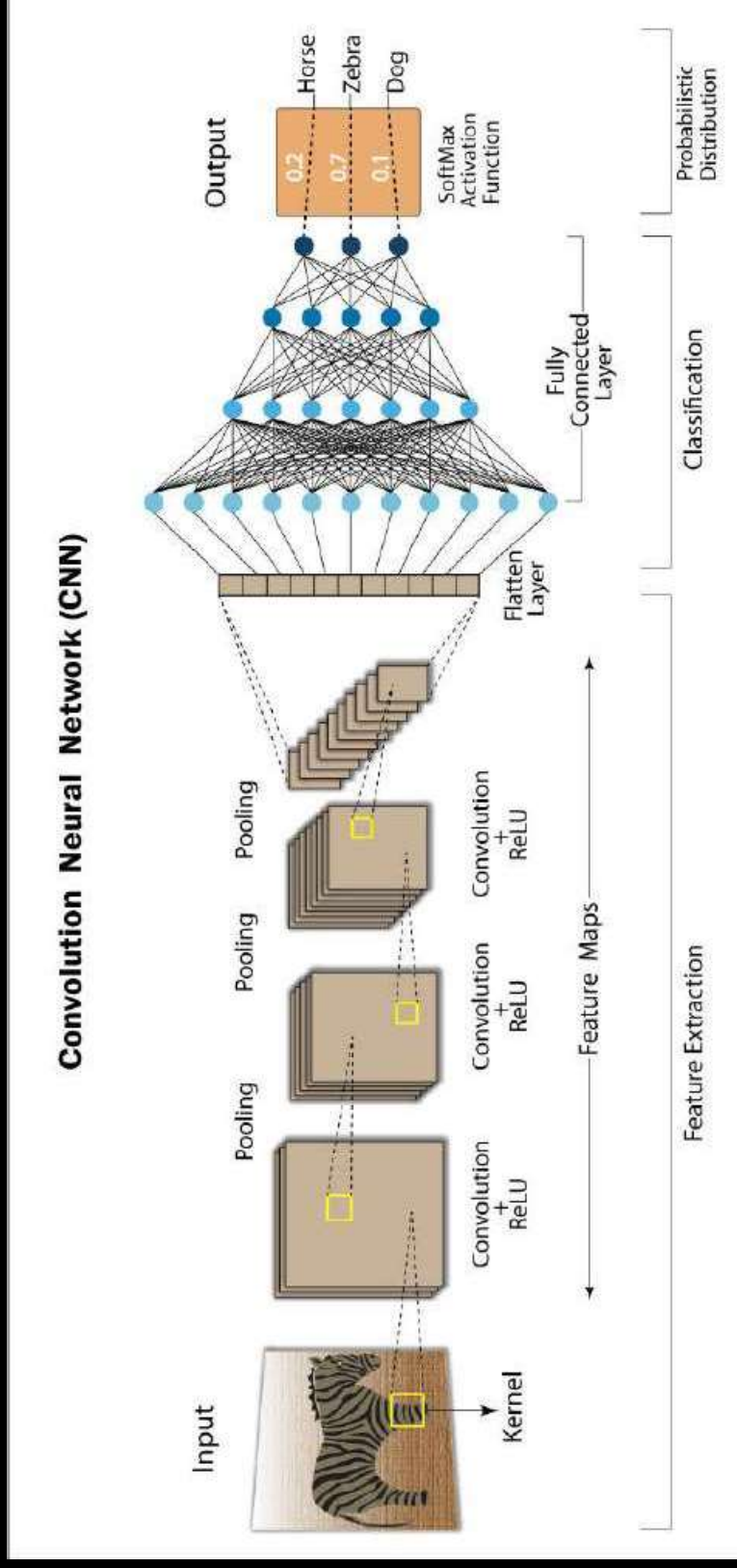
2x2 maxpooling

2	50
20	20

2) Neurone de MaxPooling :



Reconnaissance d'image



Structure type d'un réseau de classification d'images

Source : <https://www.analyticsvidhya.com/>

Résultat d'entraînement d'un réseau de neurone convolutif sur la dataset EMNIST

