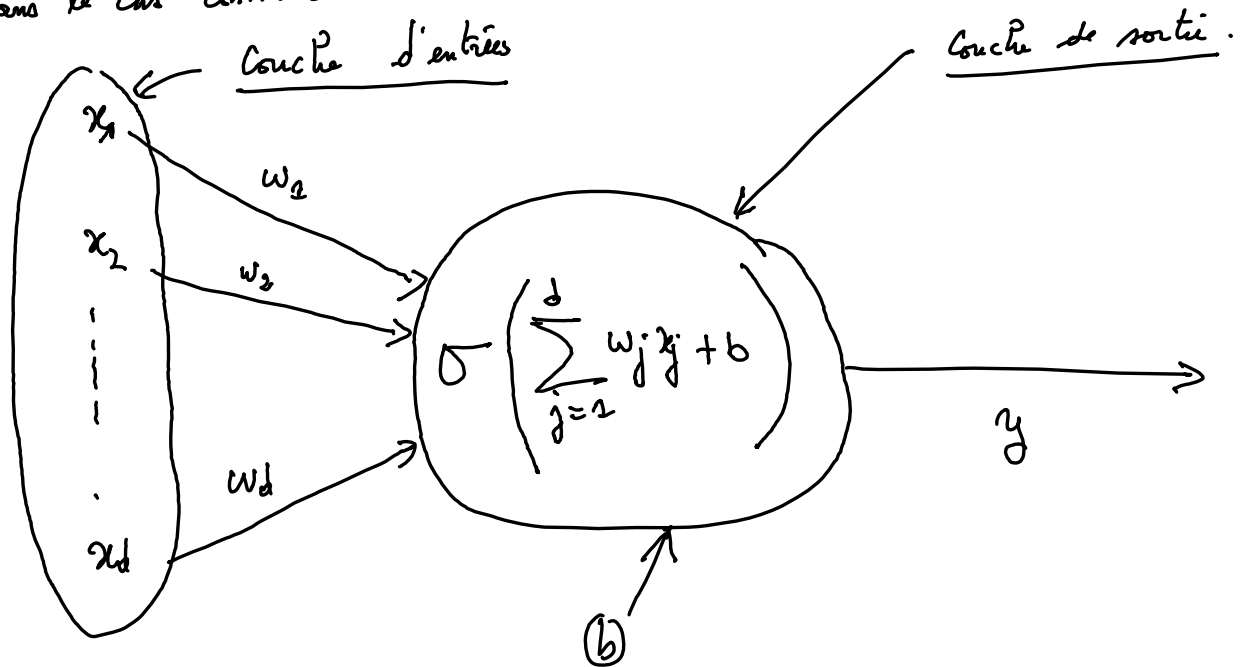
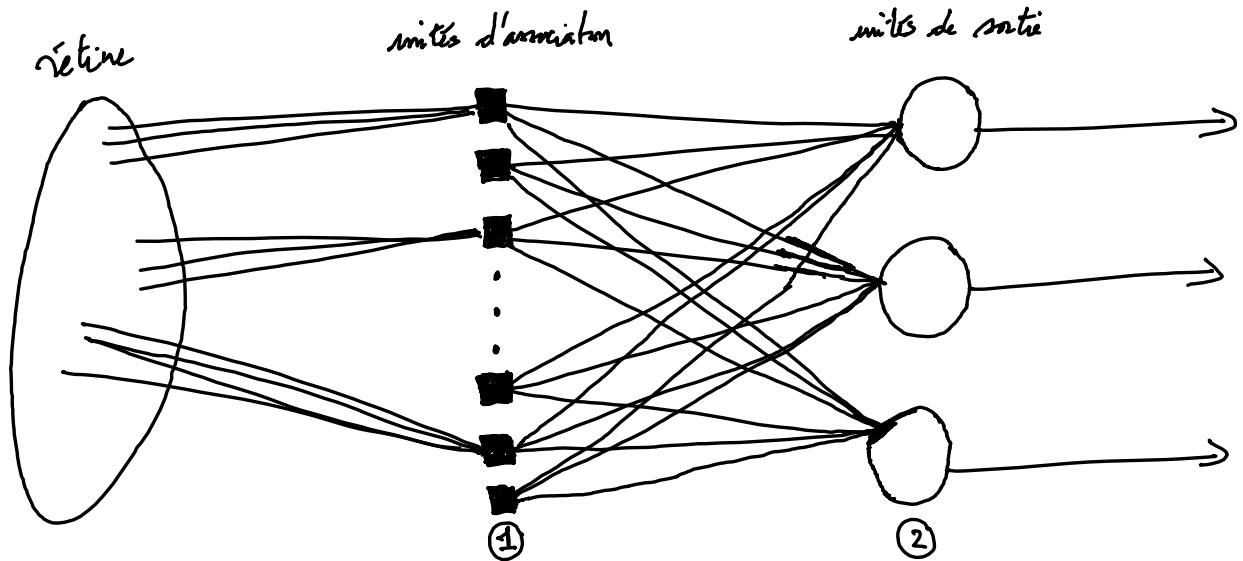


I) Architecture d'un réseau de neurones:

- Le neurophysiologiste Warren McCulloch et le mathématicien Walter Pitts ont réalisé un portrait à l'aide d'un circuit électrique.
- Un neurone de McCulloch-Pitts prend des entrées binaires, calcule une somme pondérée et renvoie 0 si le résultat est inférieur au seuil et 1 dans le cas contraire.



- À la fin des années 50, Frank Rosenblatt, psychologue à Cornell, a travaillé sur les systèmes de décision présents dans l'œil d'une mouche, qui déterminent sa réaction de fuite.
- En 1958, il a proposé l'idée d'un perceptron qu'il a appelé "Mark I Perceptron". Il s'agissait d'un système avec une relation entrée-sortie simple, modelé par un neurone de McCulloch-Pitts.



Les connexions entre ① et ② ne peuvent pas être optimisées.

Paramètres d'une perceptron:

- Fonction d'activation

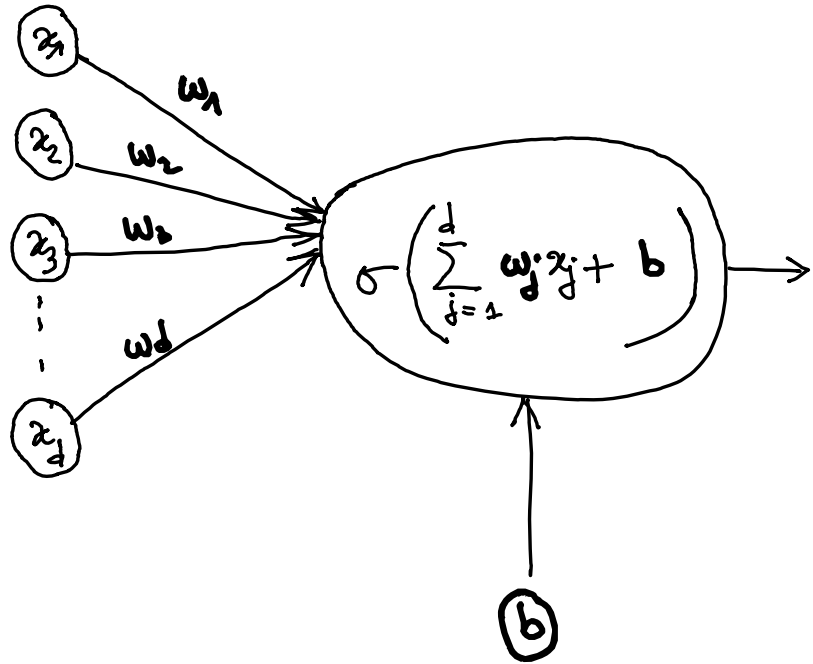
$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

• "Learnable":

- Poids $w = (w_1, \dots, w_d) \in \mathbb{R}^d$

- biais $\underline{b} \in \mathbb{R}$

Comment estimer (w, b) ?



Algorithme du Perceptron:

Nous avons des données $\mathcal{D}_n = \{(x_i, y_i), i=1, \dots, d\}$ avec $y_i \in \{-1, +1\}$.

On note $\tilde{w} = (w_1, \dots, w_d, b)$ et $\tilde{x}_i = (x_i, 1)$.

Algorithme du perceptron : (premier algorithme d'apprentissage)

* $\tilde{w} = 0$

* Pour chaque donnée (\tilde{x}_i, y_i)

- Si $y_i \langle \tilde{w}, \tilde{x}_i \rangle \leq 0$ modifier $\tilde{w} = \tilde{w} + y_i \tilde{x}_i$

- Sinon ne pas modifier \tilde{w} .

Questions:

- Pourquoi ça marche ? (intuition)

- Est-ce lié à une descente du gradient ?

Descente du gradient:

1. départ avec $\tilde{w} = 0$ (vecteur de zéros)
2. mise à jour $\tilde{w} \leftarrow \tilde{w} - \eta \nabla l(\tilde{w})$
 \downarrow
 l est la fonction de perte à minimiser.
3. arrêt si \tilde{w} ne bouge plus.

une donnée est mal-classée si $y_i \langle \tilde{w}, \tilde{x}_i \rangle \leq 0$. On veut minimiser la perte

$$l(\tilde{w}) = - \sum_{i \in \mathcal{M}_{\tilde{w}}} y_i \langle \tilde{w}, \tilde{x}_i \rangle$$

où $\mathcal{M}_{\tilde{w}}$ est l'ensemble des indices des données mal-classées par \tilde{w} .

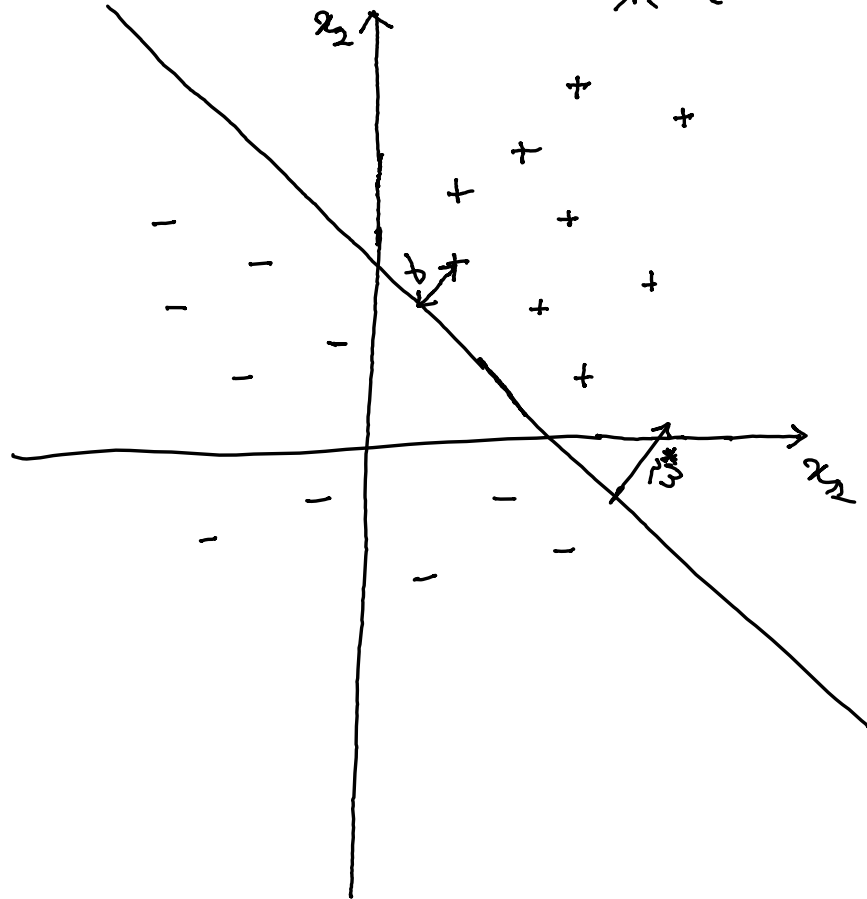
Descente du gradient: (stochastique)

1. Sélectionner au hasard un indice $i \in \mathcal{M}_{\tilde{w}}$

2. Mettre à jour $\tilde{w} \leftarrow \tilde{w} - \eta \nabla l_i(\tilde{w}) = \tilde{w} + \eta y_i \tilde{x}_i$.

Algorithme du perceptron: $\eta = 1$.

Exercice: On note $R = \max_{1 \leq i \leq n} \|x_i\|$. Soit \tilde{w}^* l'hyperplan optimal
de marge $\gamma = \min_{1 \leq i \leq n} y_i \langle \tilde{w}^*, \tilde{x}_i \rangle$ ou $\|\tilde{w}^*\| = 1$.



Théorème: (Block (1962) et Novikoff (1963))

Supposons que le jeu de données
 $D_n = \{(x_1, y_1), \dots, (x_n, y_n)\}$ est linéairement
séparable ($\gamma > 0$). On initialise
 $\tilde{w}_0 = 0$. Le nombre de mises à jour
 k de l'algorithme du perceptron
est borné par

$$k + 1 \leq \frac{1 + R^2}{\gamma^2}.$$

Inégalité de Cauchy-Schwarz: $\langle \tilde{w}^*, \tilde{w}_{k+1} \rangle \leq \|\tilde{w}_{k+1}\|_2$

Car $\|\tilde{w}^*\|_2 = 1$ avec égalité si et seulement si $\tilde{w}_{k+1} \propto \tilde{w}^*$. Par construction,

$$\begin{aligned}\langle \tilde{w}^*, \tilde{w}_{k+1} \rangle &= \langle \tilde{w}^*, \tilde{w}_k \rangle + y_i \langle \tilde{w}^*, \tilde{x}_i \rangle \\ &\geq \langle \tilde{w}^*, \tilde{w}_k \rangle + \gamma \\ &\geq \langle \tilde{w}^*, \tilde{w}_0 \rangle + (k+1)\gamma \\ &\geq (k+1)\gamma \quad \text{car } \tilde{w}_0 = 0.\end{aligned}$$

Nous avons aussi:

$$\begin{aligned}\|\tilde{w}_{k+1}\|^2 &= \|\tilde{w}_k\|^2 + \|\tilde{x}_i\|^2 + 2\langle \tilde{w}_k, y_i \tilde{x}_i \rangle \\ &\leq \|\tilde{w}_k\|^2 + R^2 + 1 \\ &\leq (k+1)(1+R^2)\end{aligned} \quad \left. \begin{array}{l} \|\tilde{x}_i\|^2 = R^2 + 1 \text{ et } \langle \tilde{w}_k, y_i \tilde{x}_i \rangle \leq 0 \\ \text{car le point } (\tilde{x}_i, y_i) \text{ est mal-classé.} \end{array} \right\}$$

$$(k+1)\gamma \leq \sqrt{(k+1)(1+R^2)} \quad \Leftrightarrow \quad k+1 \leq \frac{1+R^2}{\gamma^2}.$$

Perceptron: bilan et critiques

Bilan:

* Nous avons un jeu de données $D_n = \{(x_1, y_1), \dots, (x_n, y_n)\}$

* Nous utilisons l'algorithme du perceptron pour estimer les poids w et le terme dit de biais b .

* On va prédire en utilisant la fonction $f_{(w,b)}(x) = \mathbb{I}_{\langle w, x \rangle + b > 0}$

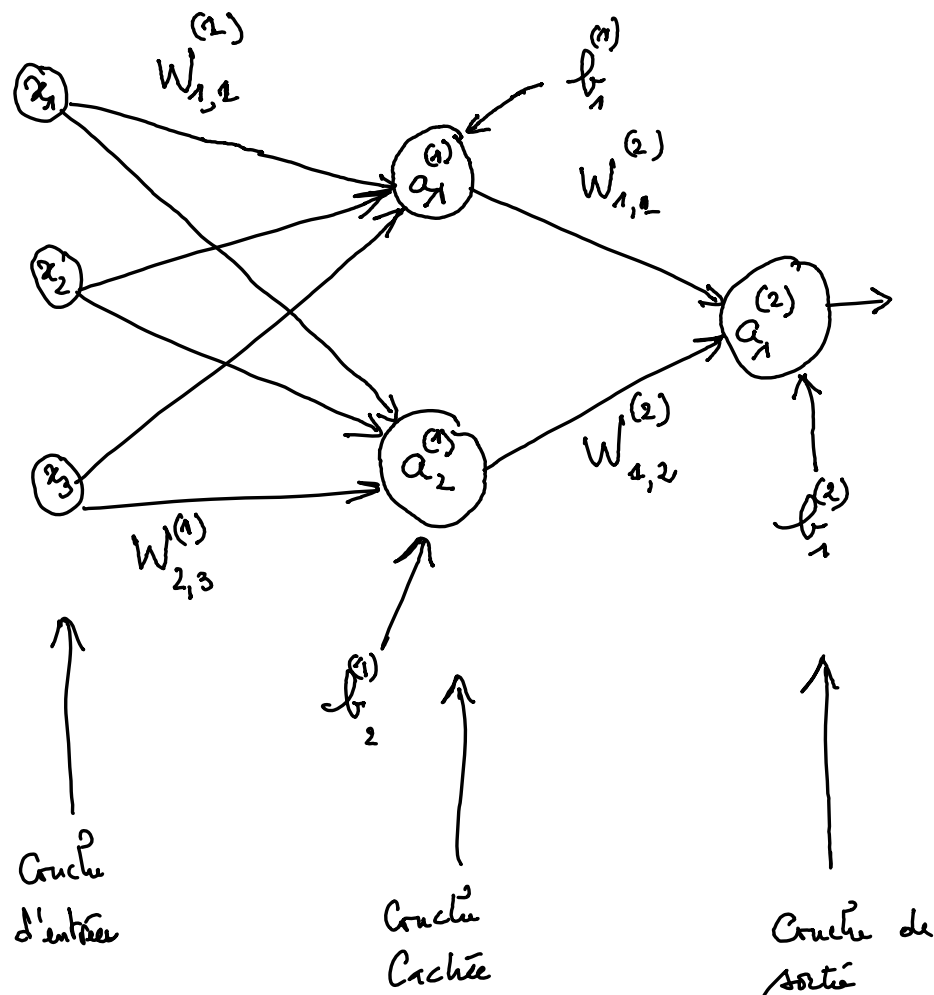
Limites:

* La frontière de décision est linéaire! (trop simple).

* L'algorithme du perceptron ne converge pas si les données ne sont pas linéairement séparables: dans ce cas, l'algorithme ne doit pas être utilisé.

* Donc, en pratique: à ne pas utiliser.

Réseau de neurones avec une couche cachée :



Notations:

1. $W_{i,j}^{(l)}$: poids entre le neurone j de la couche $l-1$ et le neurone i de la couche l .
2. $b_j^{(l)}$: biais du neurone j de la couche l .
3. $a_j^{(l)}$: sortie du neurone j de la couche l .
4. $z_j^{(l)}$: entrée du neurone j de la couche l telle que
$$a_j^{(l)} = \sigma(z_j^{(l)})$$
.

Estimation de poids et biais: (descente du gradient)

- La prédiction est donnée par $f_{\theta}(x)$
- Minimisation du risque empirique
$$\underset{\theta}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, f_{\theta}(X_i)) = \underset{\theta}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \ell_i(\theta)$$

Descente du gradient stochastique:

Tant que $\|\theta_t - \theta_{t-1}\| \geq \varepsilon$:

* On tire $I_t \subset \{1, \dots, n\}$

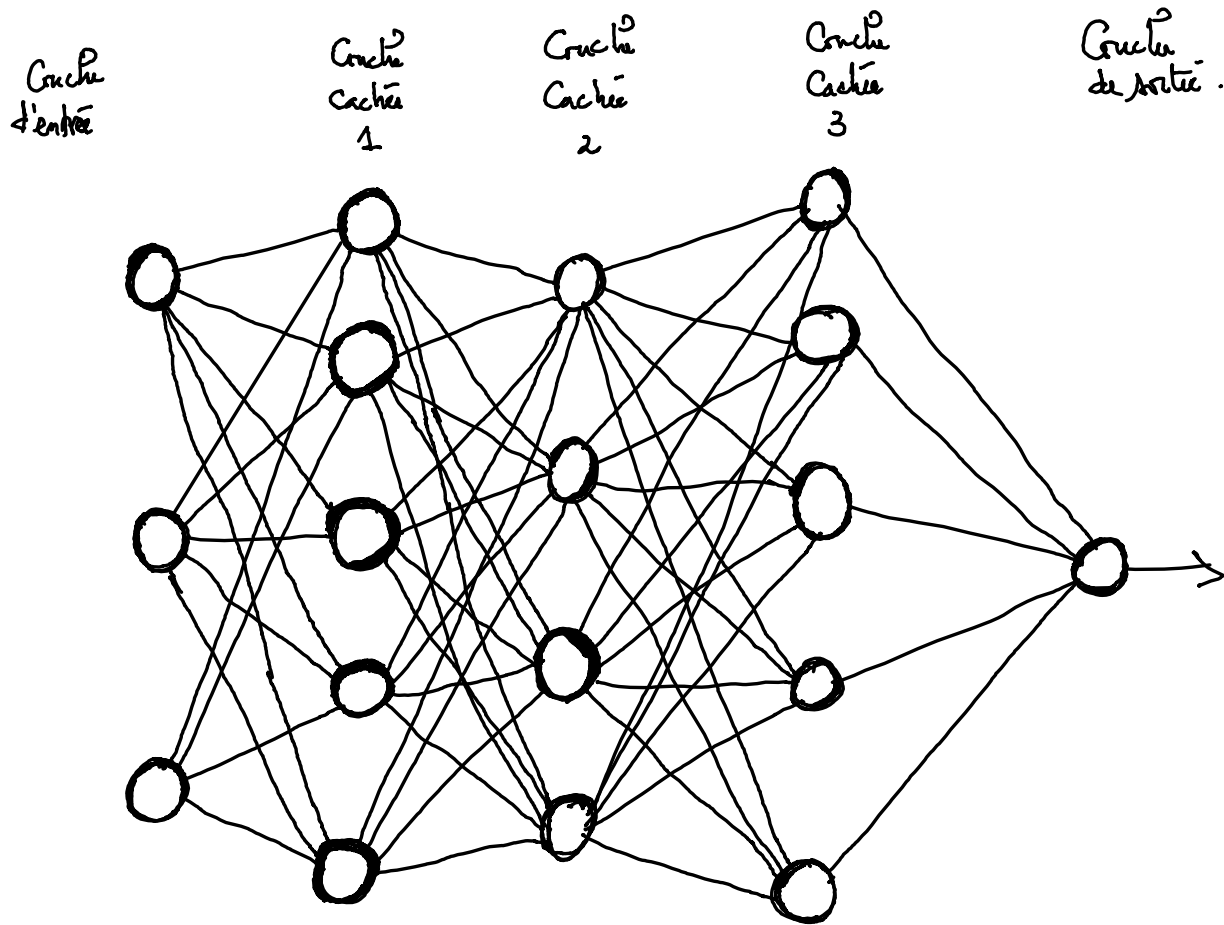
$$* \theta_{t+1} = \theta_t - \eta_t \left(\frac{1}{|I_t|} \sum_{i \in I_t} \nabla_{\theta} \ell_i(\theta_t) \right)$$

Peut-on calculer $\nabla_{\theta} \ell_i$ de manière efficace ?

Rétropropagation du gradient :

- * Connue avec les travaux de Rumelhart, McClelland, Hinton en 1986.
 - * Remonte à Werbos en 1974.
 - * C'est une formule de dérivation par chaîne.
-
- C'est la secret de fonctionnement des réseaux de neurones.
 - Au cœur du fonctionnement des réseaux de neurones profonds.

Idee de Retropropagation du gradient:



Equations de rétropropagation:

Un réseau de neurones avec L couches, avec une sortie vectorielle et une fonction de perte quadratique:

$$C = \frac{1}{2} \|y - a^{(L)}\|_2^2$$

Par définition:

$$\delta_j^{(l)} = \frac{\partial C}{\partial z_j^{(l)}}$$

Les 4 équations de la rétropropagation du gradient sont données par:

$$\delta^{(L)} = \nabla_a C \odot \sigma'(z^{(L)}) \quad \dots (1)$$

$$\delta^{(l)} = ((w^{(l+1)})^T \delta^{(l+1)}) \odot \sigma'(z^{(l)}) \quad \dots (2)$$

$$\frac{\partial C}{\partial b_j^{(l)}} = \delta_j^{(l)} \quad \dots (3)$$

$$\frac{\partial C}{\partial w_{j,k}^{(l)}} = a_k^{(l-1)} \delta_j^{(l)} \quad \dots (4)$$

Preuve: $\delta_j^{(L)} = \nabla_a C \odot \sigma'(z_j^{(L)})$, en appliquant la dérivation par chaîne :

$$\delta_j^{(L)} = \sum_k \frac{\partial C}{\partial a_k^{(L)}} \frac{\partial a_k^{(L)}}{\partial z_j^{(L)}}$$

où $z_j^{(L)}$ l'entrée du j^{ème} neurone de la couche L . Comme l'activation $a_k^{(L)}$ ne dépend seulement de l'entrée $z_j^{(L)}$, nous avons :

$$\delta_j^{(L)} = \frac{\partial C}{\partial a_j^{(L)}} \times \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}}$$

Comme $a_j^{(L)} = \sigma(z_j^{(L)})$, nous avons :

$$\delta_j^{(L)} = \frac{\partial C}{\partial a_j^{(L)}} \sigma'(z_j^{(L)})$$

C'est la j^{ème} composante de l'équation (1).

Maintenant, on veut vérifier que

$$\delta^{(l)} = [(w^{(l+1)})^T \delta^{(l+1)}] \odot \sigma'(z^{(l)}) \dots \dots \dots (2)$$

Dérivation par chaîne:

$$\begin{aligned} \delta_j^{(l)} &= \frac{\partial C}{\partial z_j^{(l)}} \\ &= \sum_k \frac{\partial C}{\partial z_k^{(l+1)}} \frac{\partial z_k^{(l+1)}}{\partial z_j^{(l)}} \\ &= \sum_k \delta_k^{(l+1)} \frac{\partial z_k^{(l+1)}}{\partial z_j^{(l)}} \end{aligned}$$

Rappelons que : $z_k^{(l+1)} = \sum_j w_{kj}^{(l+1)} \sigma(z_j^{(l)}) + b_k^{l+1}$, nous avons donc

$$\delta_j^{(l)} = \sum_k \delta_k^{(l+1)} w_{kj}^{(l+1)} \sigma'(z_j^{(l)})$$

On peut maintenant vérifier que:

$$\frac{\partial C}{\partial b_j^{(e)}} = s_j^{(e)} \dots (3)$$

Dérivation par chaîne:

$$\frac{\partial C}{\partial b_j^e} = \sum_k \frac{\partial C}{\partial z_k^{(e)}} \frac{\partial z_k^{(e)}}{\partial b_j^e}$$

$z_j^{(e)}$ dépend seulement de b_j^e et $z_j^{(e)} = \sum_k w_{jk}^{(e)} \sigma(z_k^{(e-1)}) + b_j^e$.

donc
$$\frac{\partial C}{\partial b_j^{(e)}} = s_j^{(e)}.$$

Maintenant, on s'intéresse :

$$\frac{\partial C}{\partial w_{j,k}^{(l)}} = a_k^{(l-1)} \delta_j^{(l)} \dots \quad (4)$$

Dérivation par chaîne :

$$\frac{\partial C}{\partial w_{j,k}^{(l)}} = \sum_m \frac{\partial C}{\partial z_m^{(l)}} \frac{\partial z_m^{(l)}}{\partial w_{j,k}^{(l)}}$$

Comme $z_m^{(l)} = \sum_k w_{mk}^{(l)} \sigma(z_k^{(l-1)}) + b_m^l$, nous avons

$$\frac{\partial z_m^{(l)}}{\partial w_{j,k}^{(l)}} = \sigma(z_k^{(l-1)}) \mathbb{1}_{\{m=j\}}.$$

Donc :

$$\frac{\partial C}{\partial w_{j,k}^{(l)}} = \delta_j^{(l)} \sigma(z_k^{(l-1)}) = a_k^{(l-1)} \delta_j^{(l)}.$$

Algorithme de rétropropagation:

Soit $\delta_j^{(l)} = \frac{\partial C}{\partial z_j^{(l)}}$ où $z_j^{(l)}$ est l'entrée du neurone j de la couche l .

Entraînement d'un réseau de neurones:

(a) Initialiser aléatoirement les poids et biais du réseau

(b) Pour toutes les données $(x_i)_{i \in B}$ dans le lot B .

- à répéter jusqu'à convergence.
1. Feed forward: faire passer toutes les données du lot dans le réseau et stocker les valeurs la fonction d'activation et sa dérivée de chaque neurone.
 2. Perte: Calculer la perte moyenne sur tout le lot du réseau.
 3. Rétropropagation: Calculer récursivement les vecteurs $\delta^{(l)}$ en partant de $l = L$ à $l = 1$ à l'aide des équations (1) et (2). Calculer le gradient avec les équations (3) et (4).
 4. Optimisation: Mise à jour des poids et biais à l'aide d'un pas de gradient

Vocabulaire des réseaux de neurones :

- * (Mini) batch size : nombres de données dans un feedforward/backward. Plus la taille est grande, plus vous aurez besoin de mémoire.
- * Nombre d'itérations : nombre de mixes à jour des paramètres pendant toute la durée de l'apprentissage.
- * Époque (epoch) : nombre d'itérations nécessaires pour que le réseau ait "vu" autant d'observations qu'il y en a dans le jeu de données d'apprentissage.

Par exemple : pour un jeu de données de taille 12800, si on utilise un batch size de 128, nous avons besoin de 100 itérations pour compléter "1 epoch".

Déclaration et entraînement d'un réseau de neurones :

Structure du réseau :

- * Nombre de couches / neurones par couche
- * Fonctions d'activation
- * Unité de sortie
- * Couches particulières (dropout et normalisation)

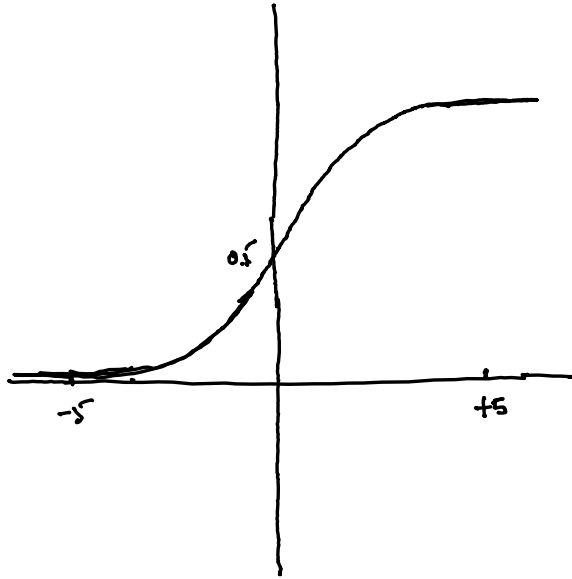
Optimisation :

- * Fonction de perte
- * Algorithme d'optimisation
- * Initialisation des poids / biais.

Nombre de neurones et couches cachées :

- * Aucune règle particulière
- * Suite des articles de recherche autour de la question et tester les propositions
- * tester des variations et regarder les performances.
- * Attention : trop de pyrométrie (sans base théorique).
- * Quelques travaux :
 - * Network pruning (Blalock et al. 2020).
 - * AgESD-Tabular (Egele et al. 2020).

Fonctions d'activation: Sigmoid



$$\sigma: x \mapsto \frac{\exp(x)}{1 + \exp(x)}$$

Commentaires:

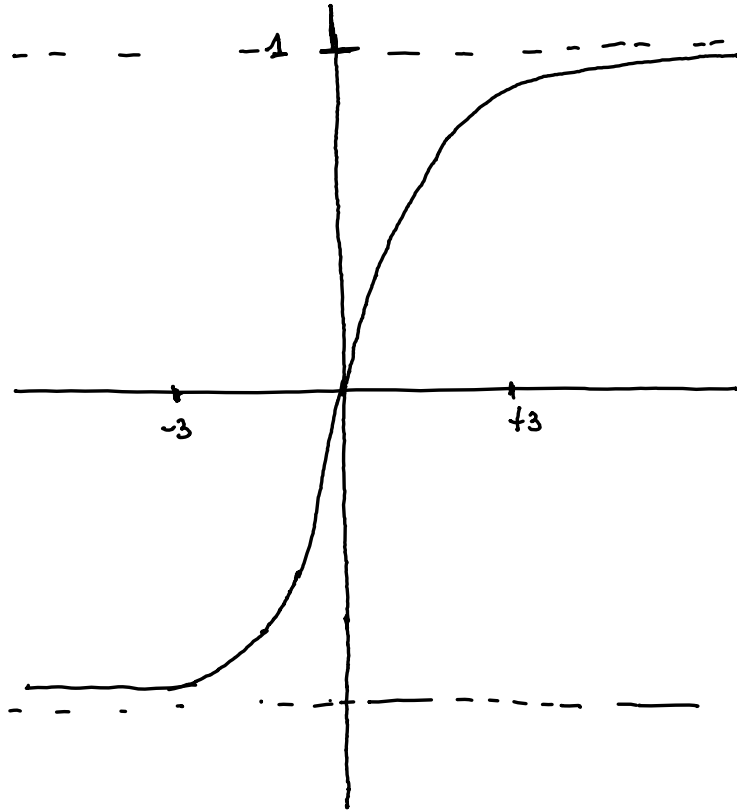
* Fonctions saturées en raison de asymptotes horizontales:

1. Gradient proche de zéro dans ces zones.

2. Normalisation des données de chaque couche pour éviter.

* $\exp(x)$ peut coûter un peu de temps.

Fonction d'activation: Tangente hyperbolique:



Commentaires:

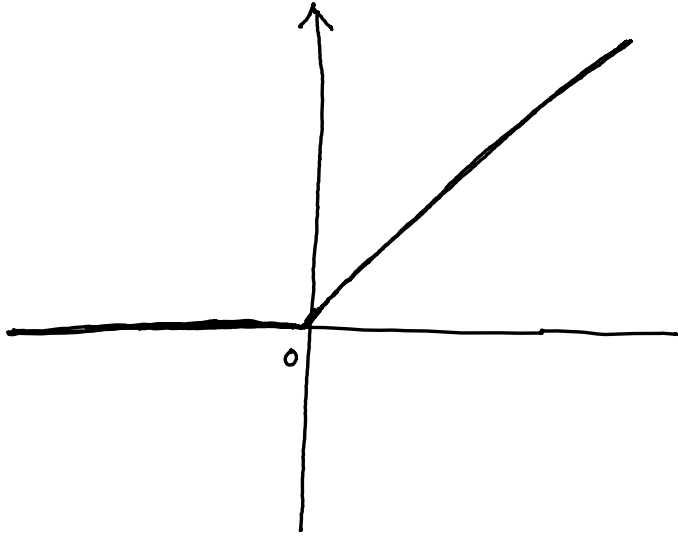
- * Fonction centrée.
- * Saturée
- * Calcul de $\exp(x)$

Notons que :

$$\tanh(x) = \sigma(2x) - 1.$$

$$\tanh: x \mapsto \frac{\exp(x) - \exp(-x)}{\exp(x) + \exp(-x)}$$

Fonction d'activation : Rectified Linear Unit (ReLU)



Commentaires:

- * Non saturée à $+\infty$

- * Saturée en $x \leq 0$.

- * Calcul pas coûteux.

$$\text{ReLU} : x \mapsto \max(0, x)$$