

Machine Learning : Réseaux de neurones

D. Cornu

Automne 2019

Le cerveau comme inspiration

« There is a fantastic existence proof that learning is possible, which is the bag of water and electricity (together with a few trace chemicals) sitting between your ears [...] wich is the squishy thing that your skull protects »

* *Stephen Marsland*

« There is a fantastic existence proof that learning is possible, which is the bag of water and electricity (together with a few trace chemicals) sitting between your ears [...] which is the squishy thing that your skull protects »

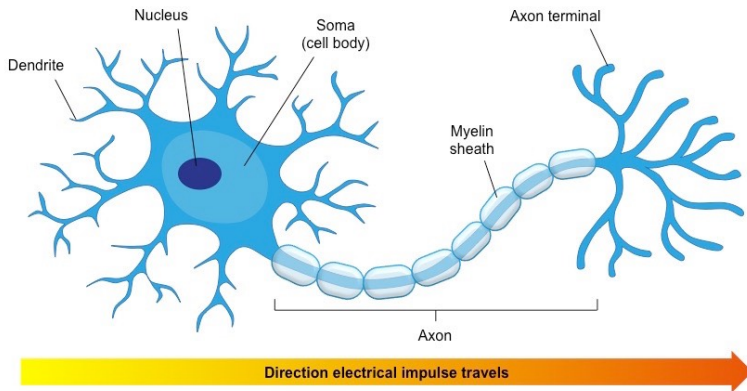
** Stephen Marsland*

Le cerveau fait exactement ce qui nous intéresse en science :

- Traiter des données bruitées ou incohérentes
- Fonctionner avec un nombre important de dimensions
- Donner un résultat le plus souvent correct
- Trouver une réponse en un temps très court
- Rester robuste malgré la perte de neurones avec l'âge

Neurone

Brique élémentaire (10^{11} dans le cerveau) = une base pour reproduire l'apprentissage



Fait la somme de signaux d'entrée. Si celle-ci est suffisante il envoie un signal le long de son axone.

Une **Synapse** est une connexion entre deux neurone (10^{14}).

Neurone = Unité de calcul simple "Tire" ou "Ne tire pas" (1 ou 0)

→ **Calculateur massivement parallèle de 10^{11} unités**

Une **Synapse** est une connexion entre deux neurone (10^{14}).
Neurone = Unité de calcul simple "Tire" ou "Ne tire pas" (1 ou 0)

→ **Calculateur massivement parallèle de 10^{11} unités**

Outils de l'apprentissage :

Les synapses représentent la **force** de la connexion entre deux neurones. Apprentissage = modification de ces liaisons
→ **plasticité**.

La **loi de Hebb** définit une règle d'apprentissage : la connexion entre deux neurones se renforce lorsqu'ils tirent au même moment
→ **conditionnement**.

Neurone Artificiel

Il est la brique élémentaire des réseaux qui seront construits.

Il est basé sur un **modèle mathématique** inspiré du neurone biologique.

C'est le modèle de **McCulloch and Pitts**.

Neurone Artificiel

Il est la brique élémentaire des réseaux qui seront construits.

Il est basé sur un **modèle mathématique** inspiré du neurone biologique.

C'est le modèle de **McCulloch and Pitts**.

Il est constitué :

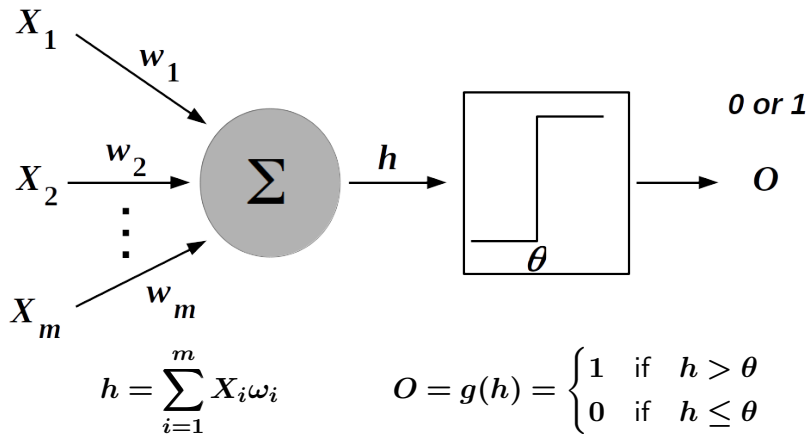
- D'un **vecteur d'entrée** X_i qui représente les différentes dimensions d'un même objet.
- D'un ensemble de **poids** w_i qui lient les entrées au neurone.
- D'une **fonction de somme** $h = \sum \dots$ qui définit comment ces poids doivent être associés aux entrées correspondantes et envoyés vers le neurone.
- D'une **fonction d'activation** $g(h)$, qui définit si le neurone doit passer dans un état actif "1" ou non "0" en fonction du résultat de la somme précédente.

Neurone Artificiel

Il est la brique élémentaire des réseaux qui seront construits.

Il est basé sur un **modèle mathématique** inspiré du neurone biologique.

C'est le modèle de **McCulloch and Pitts**.



Apprentissage supervisé → faire retrouver à un réseau la sortie attendue pour une entrée donnée.

Quel intérêt ? La sortie correcte est déjà connue à l'avance ...

→ **Généralisation** : Retrouver des "motifs" dans les données, et prédire les futures entrées dont le résultat est inconnu.

Comment changer les paramètres pour qu'un neurone apprenne ?

Apprentissage supervisé → faire retrouver à un réseau la sortie attendue pour une entrée donnée.

Quel intérêt ? La sortie correcte est déjà connue à l'avance ...

→ **Généralisation** : Retrouver des "motifs" dans les données, et prédire les futures entrées dont le résultat est inconnu.

Comment changer les paramètres pour qu'un neurone apprenne ?

Entrée et Sortie fixées, l'apprentissage repose donc uniquement sur les poids W_i et la limite d'activation θ

Apprentissage

Il s'agit d'un modèle **Supervisé**, chaque vecteur d'entrée est donc associé à une réponse attendue : la **target t** .

Supposons que l'on présente un vecteur d'entrée au neurone et que celui-ci ne donne pas la réponse attendue.

Il y a **m poids w_i** (avec i de 1 à m) qui sont connectés à ce neurones, correspondant à chaque noeud d'entrée.

Comment modifier les poids ? :

- Si le neurone tire alors qu'il ne devrait pas, les poids sont trop grands
- Si le neurone ne tire pas alors qu'il devrait, les poids sont trop petits

On calcule une erreur : **$y_k - t_k$**

La différence entre la sortie et la cible pour le neurone k

Cette différence est multipliée par la valeur de chaque entrée et utilisée pour **changer le poid associé** **$\Delta w_{ij} = -\eta (y_k - t_k) \times x_i$**

Taux d'apprentissage

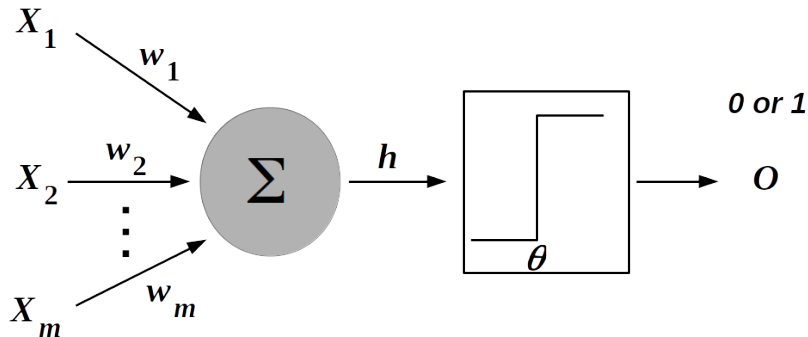
On définit un **taux d'apprentissage** η qui permet de quantifier la vitesse à laquelle l'algorithme va modifier les poids.

- Une valeur trop importante de ce taux rends l'algorithme instable.
- Une valeur trop faible ralentie l'apprentissage.

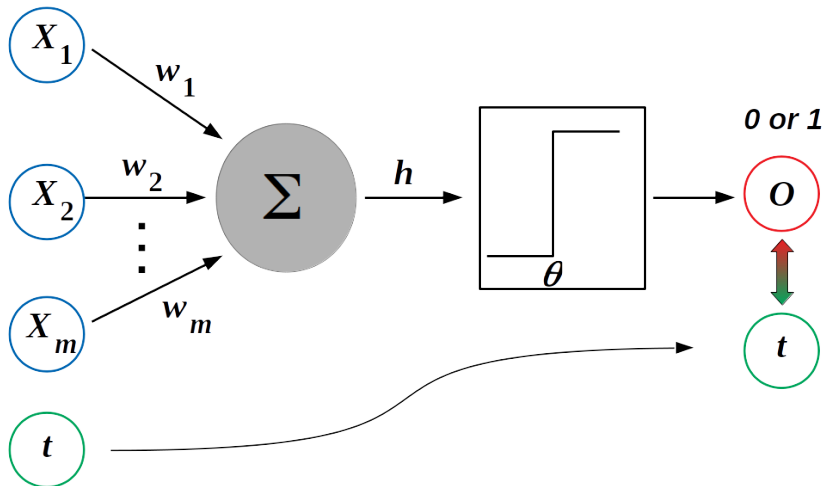
Une valeur typique se situe entre $0.1 < \eta < 0.4$ ce qui permet de rendre l'algorithme plus stable mais aussi plus résistant aux erreurs (bruit).

Le choix réside dans la qualité attendue des données d'apprentissage.

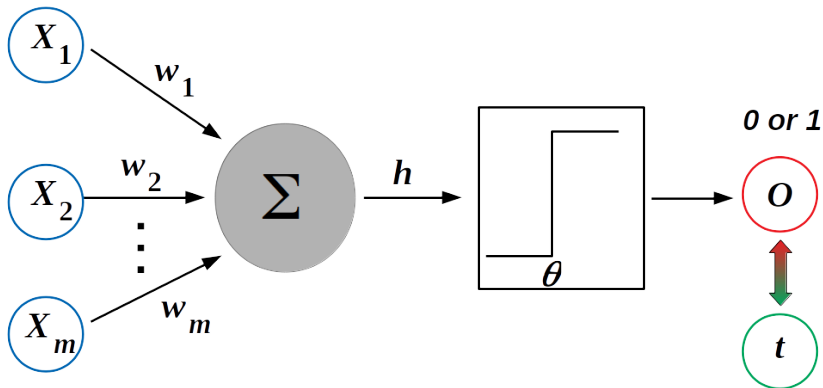
Apprentissage illustré



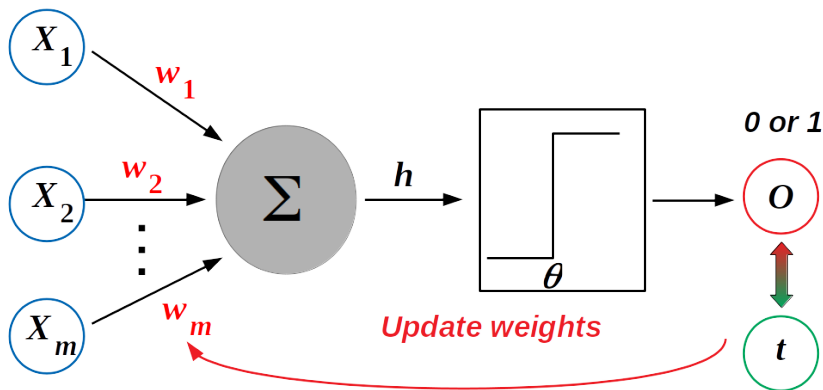
Apprentissage illustré



Apprentissage illustré



Apprentissage illustré

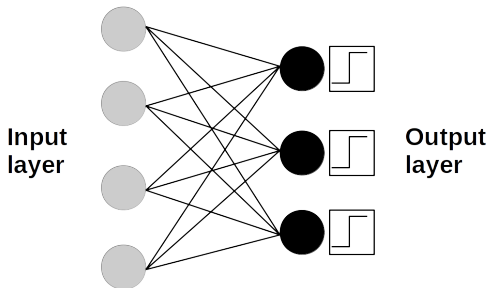


$$\omega_i \leftarrow \omega_i - \eta (o - t) \times X_i$$

Le Perceptron

Un seul neurone est dans la majorité des cas insuffisant pour représenter correctement un problème. La façon la plus simple d'en ajouter est de le faire sur une "couche/layer". Ces neurones restent indépendants, chacun représentant une **séparation linéaire** dans l'espace des paramètres d'entrée, mais l'ensemble peut servir à encoder une information.

- Les neurones sont indépendants entre eux
- Les poids sont séparés pour chaque neurone
- Le nombre de dimensions en entrée et le nombre de neurones (sorties) sont indépendants
- Les dimensions d'entrée et de sortie sont imposées par les données
- La sortie est un schéma de 0 et de 1 qui encode l'information voulue



Vecteur d'entrée connecté de manière pondérée à des neurones de McCulloch

The Perceptron

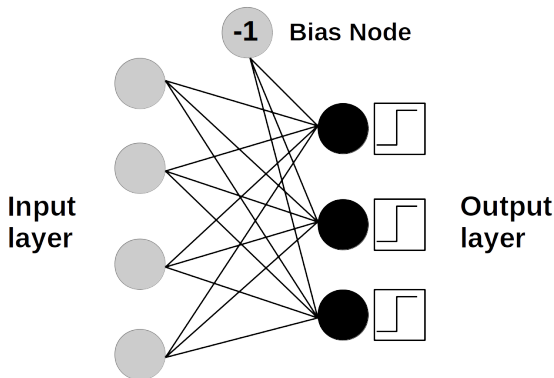
L'entrée biaisée

Problème :

Si toutes les entrées sont à 0 alors tous les neurones en sortie auront **le même comportement...** Pour contrer cet effet la limite d'activation peut être ajustée, mais difficile à mettre en œuvre

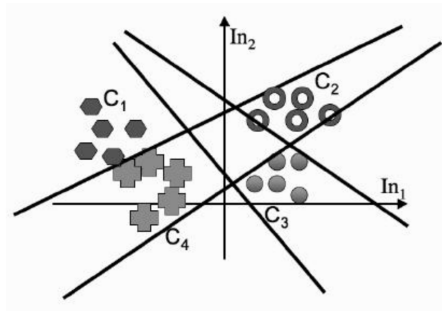
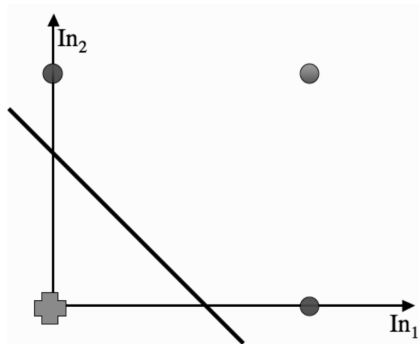
Solution :

Ajouter un **noeud d'entrée** qui a une valeur fixée à -1 et des poids associés. Définira le comportement pour des entrées proches de zéro.



Note sur la séparabilité linéaire

Les neurones du perceptron tel que décrits ici ne sont capables que de tracer des séparations linéaires entre les classes étudiées.



Il est donc nécessaire de représenter les données en entrée de manière astucieuse en ajoutant un nombre de dimensions suffisant afin de garantir leur séparabilité.

Un neurone réel reste assez différent :

- La somme des inputs peut etre non linéaire
- La sortie n'est pas binaire mais est une séquence d'impulsion ce qui contient de l'information supplémentaire
- La limite de tir varie au cours du temps
- L'interrogation des neurones n'est pas séquentielle (mise à jour Asynchrone)
- Les poids peuvent être négatifs ou positifs mais pas de passage de l'un à l'autre
- Les synapses peuvent revenir sur le neurone d'origine (feedback)
- Après avoir tiré, un neurone à un temps de recharge

Le model théorique est cependant suffisant pour apprendre des images, représenter des fonctions, faire du classement, ...

Perceptron : Algorithme

- Initialisation :
 - Définir tous les poids w_i avec de petites valeurs aléatoire (positives et négatives)
- Entraînement
 - pour T itérations ou jusqu'à ce que la sortie soit correcte
 - ★ Pour chaque vecteur d'entrée
 - Calculer l'activation de chaque neurone j avec la fonction d'activation g :

$$y_j = g \left(\sum_{i=0}^m w_{ij} x_i \right) = \begin{cases} 1 & \text{if } \sum_{i=0}^m w_{ij} x_i > 0 \\ 0 & \text{if } \sum_{i=0}^m w_{ij} x_i \leq 0 \end{cases} \quad (1)$$

- Mettre à jour les poids individuellement avec :

$$w_{ij} \leftarrow w_{ij} - \eta (y_j - t_j) \cdot x_i \quad (2)$$

- Rappel
 - Calculer l'activation de chaque neurone j pour chaque vecteur

Première application : la porte OU

Un exemple très simple de l'application de cet algorithme est de tenter de lui apprendre **la porte logique "OU"**.

Le perceptron prend donc une **entrée à deux dimensions** et ne possède qu'**un seul neurone** qui doit donner la sortie attendue sous la forme d'un 0 ou d'un 1.

Codez donc un premier programme qui declare un tableau contenant les entrées possibles pour cette porte logique, et un tableau pour les cibles correspondantes.

Declarez les variables necessaire à l'algorithme, initialisez les deux poids à des petites valeurs aléatoires, et tentez de faire un entrainement sur quelques itérations. **N'oubliez pas le noeud de biais !**

Affichez la sortie, à chaque iteration pour observer la convergence.

Correction porte OU

```
1 program logic_or
2   implicit none
3
4   integer, dimension(4,3) :: input
5   integer, dimension(4) :: targ
6   integer :: output
7   real :: h, learn_rate = 0.1
8   real, dimension(3) :: weights
9   integer :: i, j, t
10
11  input = reshape((/0, 0, 1, 1, 0, 1, 0, 1, -1,
12                 -1, -1, -1 /), shape(input))
13  targ = (/0, 1, 1, 1/)
14  call random_number(weights)
15  weights(:) = weights(:) * (0.02) - 0.01
16
17  ! #####
18  ! #####
19  ! ##### Main training loop #####
20  do t = 1, 5
21    write(*,*) "Iteration :", t
22    ! #####
23    ! Testing the result of the network with a
      forward
24    ! #####
25    ! #####
26    do i=1, 4
27      !Forward phase
28      h = 0.0
29      do j=1, 3
30        h = h + weights(j)*input(i,j)
31      end do
32
33      if(h > 0) then
34        output = 1
35      else
36        output = 0
37      end if
38
39      write(*,*) "Input :", input(i,:)
40      write(*,*) "Target :", targ(i)
41      write(*,*) "Output :", output
42      write(*,*) " "
43    end do
44
45    ! #####
46    ! ##### Training on all data once #####
47    ! #####
48    ! #####
49    do i=1, 4
50      !Forward Step
51      h = 0.0
52      do j=1, 3
53        h = h + weights(j)*input(i,j)
54      end do
55
56      if(h > 0) then
57        output = 1
58      else
59        output = 0
60      end if
61
62      !Back-propagation phase
63      do j=1, 3
64        weights(j) = weights(j) - learn_rate*(
65          output - targ(i))*input(i,j)
66      end do
67    end do
68  end program logic_or
```

The Pima Indian Dataset

Fait parti du UCI Machine Learning repository, très utile pour récupérer des données de test.

Donne 8 mesures effectuées sur une population de natifs américains nommés les Pimas, et les classe selon que la personne soit atteinte de diabète ou non.

Malgré ses limitations le perceptron une fois correctement entraîné devrait être capable de trouver le résultat attendu dans 60 à 70% des cas.

Modifiez votre programme précédent pour lire les données pima. Adaptez le nombre d'entrées et les calculs effectués. Mesurez la "précision" de votre algorithme et son évolution au fil des itérations.

Correction pima

```
1 program logic_or
2   implicit none
3
4   integer, parameter :: nb_dat = 768
5   real, dimension(nb_dat,9) :: input
6   integer, dimension(nb_dat) :: targ
7   integer :: output, precis
8   real :: h, learn_rate = 0.1
9   real, dimension(9) :: weights
10  integer :: i, j, t
11
12
13  open(10, file="pima-indians-diabetes.data")
14
15  do i=1, nb_dat
16    read(10,*) input(i,1:8), targ(i)
17  end do
18  input(:,9) = -1.0
19
20
21  call random_number(weights)
22  weights(:) = weights(:) * (0.02) - 0.01
23
24  ! #####
25  !           Main training loop
26  ! #####
27  do t = 1, 50
28    write(*,*) "Iteration :", t
29    ! #####
30    ! Testing the result of the network with a
31    ! forward
32    ! #####
33    precis = 0.0
34    do i=1, nb_dat
35      !Forward phase
36      h = 0.0
37      do j=1, 9
38        h = h + weights(j)*input(i,j)
39      end do
40      if(h > 0) then
41        output = 1
42      else
43        output = 0
44      end if
45      if(output == targ(i)) then
46        precis = precis + 1
47      end if
48    end do
49    write(*,*) "Precision rate :", real(precis)/
50      real(nb_dat)
51
52    ! #####
53    !           Training on all data once
54    ! #####
55    ! #####
56    do i=1, nb_dat
57      !Forward phase
58      h = 0.0
59      do j=1, 9
60        h = h + weights(j)*input(i,j)
61      end do
62
63      if(h > 0) then
64        output = 1
65      else
66        output = 0
67      end if
68
69      !Back-propagation phase
70      do j=1, 9
71        weights(j) = weights(j) - learn_rate*(
72          output-targ(i))*input(i,j)
73      end do
74    end do
75  end do
76 end program logic_or
```

Amélioration : Préparer les données

Un neurone donne une sortie 0 ou 1 il est donc assez évident de transformer la sortie attendue pour qu'elle puisse être encodée avec ce type de sortie.

Il est également utile de faire quelques changements sur les entrées tel que celles ci se trouvent entre -1 et 1 . On parle de **normalisation** des données.

Il est aussi possible de faire des groupes dans certains paramètres, par exemple l'âge, en faisant des catégories 21-30,31-40,... Plus facilement reconnaissable et plus logique.

Autre modification, **la réduction dimensionnelle**, qui consiste à retirer successivement certains paramètres en entrée et de voir si le résultat change ou non.

Correction pima optimisation

```
1
2
3  do i = 0, nb_dat
4      if(input(i,1) > 8) input(i,1) = 8
5      input(i,8) = int(mod(input(i,8) - 30, 10.0))
6      if(input(i,8) > 5) input(i,8) = 5
7  end do
8
9  do i=1, 8
10     input(:,i) = input(:,i) - (sum(input(:,i))/nb_dat)
11     input(:,i) = input(:,i)/maxval(abs(input(:,i)))
12 end do
```

Ce jeu de données vient également de l'UCI repository. Il donne 4 mesures effectuées sur des iris et les classe en 3 catégories.

La meilleure façon pour faire une classification qui n'est pas binaire est d'avoir autant de neurone que de classes de sortie. La target voulue sera alors encodée de telle manière que le neurone associé à la classe vaudra 1 et que les autres vaudront 0.

Exemple : classe 1 : 1 0 0 - classe 2 : 0 1 0 - classe 3 : 0 0 1

Modifiez votre algorithme pour qu'il lise ce nouveau jeu de données et associe correctement les cibles. Adaptez le cœur de l'algorithme comme expliqué plus haut pour prendre en charge plusieurs neurones de sortie indépendants.

Correction Iris dataset

Voir code (trop long pour être affiché)

Note sur le sur-apprentissage

Entraîner **trop longtemps** l'algorithme peut avoir pour conséquence un apprentissage du bruit contenu dans les données. L'algorithme perd alors trace de la tendance...

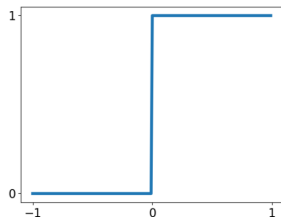
Soit limiter l'entraînement manuellement, soit surveiller l'évolution de l'apprentissage.

De plus il est important de **mélanger** les données (en conservant la bonne association de cibles) à chaque iteration sur l'ensemble du dataset, afin d'éviter un biais dans l'organisation des données.

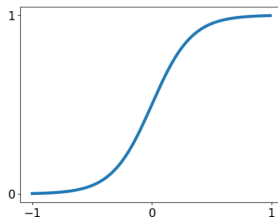
Dans le cas du perceptron ces précautions ne sont pas vraiment nécessaires puisque la convergence de l'algorithme est dominée par les effets non linéaires manquants.

Nécessite deux amélioration principales :

- Changement de la fonction d'activation : Passer à une fonction **continue dérivable** non binaire. Bon exemple : **la Sigmoid**.



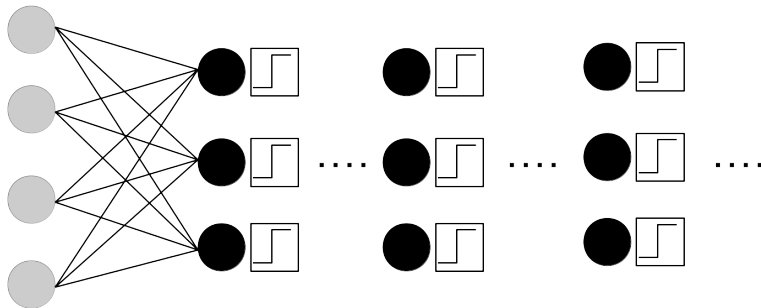
Threshold



Sigmoid

Nécessite deux amélioration principales :

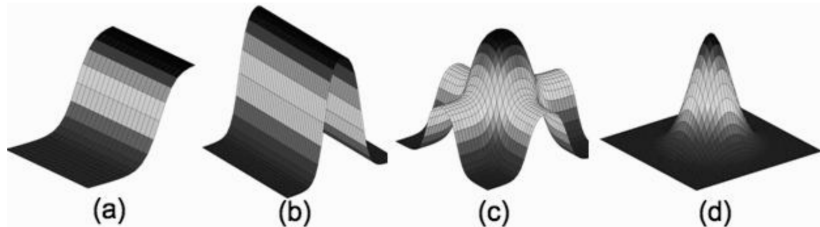
- Changement de la fonction d'activation : Passer à une fonction **continue dérivable** non binaire. Bon exemple : **la Sigmoid**.
- Ajouter des neurones sur une autre couche pour permettre les combinaisons **non linéaires**.



Nécessite deux améliorations principales :

- Changement de la fonction d'activation : Passer à une fonction **continue dérivable** non binaire. Bon exemple : **la Sigmoid**.
- Ajouter des neurones sur une autre couche pour permettre les combinaisons **non linéaires**.

Ces deux améliorations font de cet algorithme un **"Universal Function Approximator"**.



Multi Layer Perceptron

Ce nouvel algorithme, le MLP, ajoute les neurones sous la forme de couches dites "cachées/hidden"

Problème : Comment estimer l'erreur sur la couche cachée ?

Multi Layer Perceptron algorithm

$$\frac{\partial E}{\partial \omega_{ij}} = \frac{\partial E}{\partial h_j} \frac{\partial h_j}{\partial \omega_{ij}} \quad \delta_l(j) \equiv \frac{\partial E}{\partial h_j} = \frac{\partial E}{\partial a_j} \frac{\partial a_j}{\partial h_j} = g'(a_j) \sum_k \omega_{kj} \delta_{l-1}(k)$$

Les changements précédents se traduisent par les modifications suivantes dans l'algorithme du perceptron :

- Changement de la fonction d'activation de tous les neurones tel que :

$$a_k = g(h_k) = 1 / (1 + \exp(-\beta h_k))$$

avec β un réel positif supérieur ou égal à 1.0 définissant la "pente".

- Ajout d'une couche cachée, avec son noeud de biais associé.
- Changement de la "back propagation" avec $E(a, t) = \frac{1}{2} \sum_{k=1}^N (a_k - t_k)^2$:

$$\delta_o(k) = \beta(a_k - t_k)a_k(1 - a_k)$$

$$\delta_h(j) = \beta a_j(1 - a_j) \sum_{k=1}^N \delta_o(k) \omega_{jk}$$

$$\omega_{jk} \leftarrow \omega_{jk} - \eta \delta_o(k) a_j$$

$$v_{ij} \leftarrow v_{ij} - \eta \delta_h(j) x_i$$

Remarques :

Les equations précédentes sont valable pour le cas ou tout les neurones ont une fonction d'activation sigmoïd et pour une erreur quadratique classique.

Dans le cas général les neurones cachés seront toujours sous la forme de sigmoid ou d'une fonction proche. Néanmoins les neurones de sortie peuvent être adaptés pour etre plus performants sur certaines taches.

Cet algorithme est capable de faire de la régression, pour cela il faut changer la fonction d'activation du/des neurones de sorties vers une fonction linéaire $o = g(h) = h$. Il faut donc adapter le calcul de $\delta^o = (o - t)$.

Pour de la classification, une sortie ou chaque neurone correspond à une classe est plus adapté mais on voudrait normaliser cette sortie. Pour cela on peut utiliser la fonction d'activation SOFT-MAX telle que : $a_k = g(h_k) = \frac{\exp(h_k)}{\sum_{k=1}^N \exp(h_k)}$. Il faut alors adapter δ^o , qui se simplifie par la même fonction que pour la sigmoid mais sans le parametre β .

Sur apprentissage : Séparation des données

Apprentissage = voir plusieurs fois l'échantillon

Entraîner **trop longtemps** cause un **surapprentissage** !

→ apprend du bruit, ou des spécificités de l'échantillon.

L'algorithme perd alors trace de la tendance... **Pour contrer ce problème il convient de séparer les données :**

- **Entrainement** : Majorité des données, sert à la phase d'entraînement de l'algorithme.
- **Validation** : partie plus modeste des données, sert à mesurer l'erreur au cours de l'apprentissage. Si elle augmente de manière significative il faut stopper l'apprentissage.
- **Test** : permet de contrôler sur des données indépendantes la qualité de l'apprentissage une fois celui-ci terminé.

Des proportions usuelles sont 50 : 25 : 25 ou encore 60 : 20 : 20