

Réseaux de Neurones Artificiels

M1 Miage 2016–2017 *Intelligence Artificielle*

Stéphane Airiau



Ce cours suit le chapitre 18 de “Artificial Intelligence : A modern Approach”.

-
- S'inspirer du cerveau humain pour construire un système
 - Ici, on s'inspire des neurones
 - ↳ bâtir des réseaux de neurones artificiels
 - vieille idée, McCulloch et Pitts ont présenté un modèle mathématique d'un neurone dès 1943.

- le cerveau : réseau interconnecté très complexe de neurones
- Réseau de Neurones Artificiels (RNA) : réseau densément connecté de simples unités de calcul
- le cerveau
 - $\approx 10^{11}$ neurones
 - chaque neurone est connecté à 10,000 autres neurones
 - le temps de réaction d'un neurone est autour de 10^{-3} s (lent par rapport à un ordinateur)
 - 0.1s pour réaliser une tâche de reconnaissance (ex : montre la photo d'une personnalité ou d'un ami)
 - au plus quelques centaines d'étapes pour le calcul.
 - calcul vraisemblablement massivement parallèle
- réseaux de neurones pour comprendre mieux le cerveau humain
- réseaux de neurones comme source d'inspiration pour un mécanisme d'apprentissage efficace

- domaine de recherche ancien
- ex dès 1993, un RNA a été utilisé pour gérer le volant d'une voiture roulant sur une autoroute (aux États Unis)
 - input : images d'une caméra avec une résolution de 30x32 pixels ↪ chaque pixel est connecté à un neurone d'entrée
 - 4 neurones cachés
 - 30 neurones pour la sortie (un neurone va correspondre à un angle du volant)

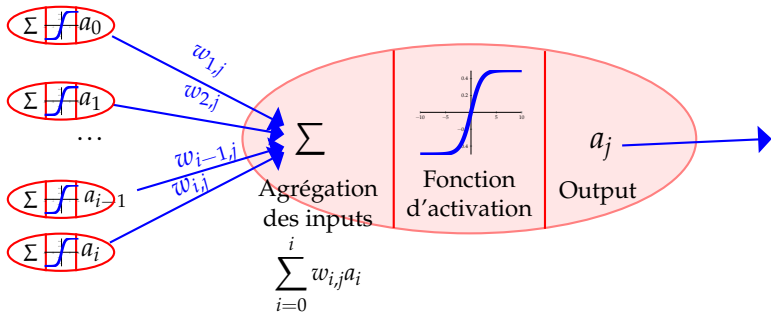
Situations pour utiliser un RNA

- entrée est un vecteur de large dimension de valeurs discrètes ou continues (e.g. sortie d'un capteur)
- la sortie est discrète ou continue
- la sortie est un vecteur de valeurs
- les valeurs d'entrées peuvent avoir du bruit (ou contenir des erreurs)
- on n'a pas d'hypothèse sur la fonction à apprendre
- il n'y a pas besoin qu'un humain puisse lire le résultat (par ex pour le vérifier ou comprendre ce qui est appris)
- un temps d'apprentissage long est acceptable
- l'utilisation du RNA est rapide

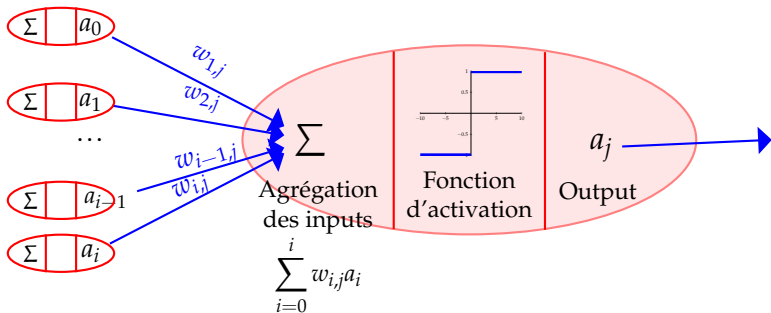
Exemples :

- apprendre à reconnaître des sons
- classification d'images
- prédictions financières

Modèle d'un neurone : le perceptron



Modèle d'un neurone : le perceptron



- perceptron représente un hyperplan $\vec{w} \cdot \vec{x} = 0$ et les instances d'un côté de l'hyperplan seront positives et négatives de l'autre
- on peut représenter des fonction booléennes avec un seul perceptron, mais pas toutes (XOR).
- avec un réseau de perceptrons, on peut représenter toutes les fonctions booléennes

Apprendre pour un perceptron

- on doit apprendre les poids w_i entre les entrées $1, 2, \dots, n$ et le perceptron.
- on peut commencer par des poids pris au hasard
- on présente une série d'instances et on corrige les poids après chaque erreur
 - t est la vraie classification de la sortie
 - o est la sortie générée par le perceptron
 - x_i est l'état de l'entrée numéro i

$$w_i \leftarrow w_i + \Delta w_i$$

où $\Delta w_i = \eta(t - o)x_i$

- η est le taux d'apprentissage
(une petite valeur comme 0.1 est habituellement utilisée)
- si les données peuvent être apprises par un perceptron (on dit que les données sont alors linéairement séparables), ce simple algorithme **converge**.

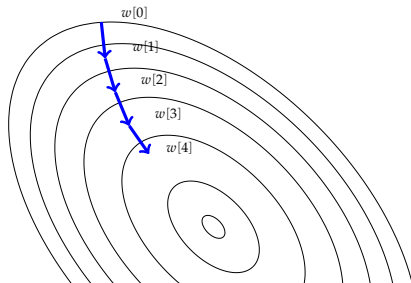
- Si les données sont linéairement séparables, la méthode précédente converge
- et sinon ? La règle suivante va garantir convergence à la meilleure approximation
- on va utiliser le principe la **descente de gradient** (la plus grande pente)
- L'erreur d'apprentissage est souvent mesurée par

$$E(\vec{w}) = \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

où

- D est l'ensemble de exemples d'apprentissage
- t_d est la vraie classification de l'instance d
- o_d est la réponse du perceptron pour l'instance d

Descente de gradient



le gradient indique la direction de la pente la plus forte

La règle de mise à jour sera donc :

$$\vec{w} \leftarrow \vec{w} + \Delta \vec{w}$$

$$\text{où } \Delta \vec{w} = -\eta \nabla E(\vec{w})$$

$$w_i \leftarrow w_i + \Delta w_i$$

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i}$$

- η est le taux d'apprentissage qui détermine la taille du pas que l'on fait en descendant
- le signe négatif indique que l'on veut descendre

Heureusement, calculer la dérivée est facile ici !

$$\begin{aligned}\frac{\partial E}{\partial w_i} &= \frac{\partial}{\partial w_i} \left(\frac{1}{2} \sum_{d \in D} (t_d - o_d)^2 \right) \\ &= \frac{1}{2} \sum_{d \in D} \frac{\partial}{\partial w_i} (t_d - o_d)^2 \\ &= \frac{1}{2} \sum_{d \in D} 2(t_d - o_d) \frac{\partial}{\partial w_i} (t_d - o_d) \\ &= \sum_{d \in D} (t_d - o_d) \frac{\partial}{\partial w_i} (t_d - \vec{w} \cdot \vec{x}_d) \\ \frac{\partial E}{\partial w_i} &= \sum_{d \in D} (t_d - o_d) (-x_{i,d})\end{aligned}$$

La règle de mise à jour est donc

$$\Delta w_i = \eta \sum_{d \in D} (t_d - o_d) x_{i,d}$$

Algorithme de descente de gradient

algorithme pour deux niveau d'unité avec des fonctions sigmoids.

```
1 initialise chaque  $w_i$  avec une valeur au hasard
2 tant que l'erreur est trop grande
3   Pour chaque  $i \in \{1, \dots, n\}$ 
4      $\Delta w_i = 0$ 
5     Pour chaque exemple  $(\vec{x}, t) \in D$ 
6       calculer la sortie  $o$ 
7       Pour chaque  $i \in \{1, \dots, n\}$ 
8          $\Delta w_i = \Delta w_i + \eta(t - o)x_i$ 
9       Pour chaque  $i \in \{1, \dots, n\}$ 
10       $w_i \leftarrow w_i + \Delta w_i$ 
```

Descente de gradient stochastique

- la descente peut être lente
- s'il y a plusieurs minimaux locaux, on n'a pas de garantie de trouver le minimum global

On utilise souvent une variante qui consiste à mettre à jour les poids après l'examen de chaque point (et pas après de les avoir tous examinés)

- c'est souvent une bonne approximation
- demande un peu moins de calculs
- permet parfois d'éviter de tomber dans un minimum local ➡ plus de chance de tomber dans le minimum global

Il faut passer au réseau de neurones!

- on veut représenter des fonctions non linéaires
- avec la discontinuité du perceptron, on ne peut pas calculer de dérivées
- on remplace la fonction d'activation par une fonction **sigmoid** (ou fonction logistique) qui est une approximation continue et différentiable de la fonction seuil

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

L'erreur d'apprentissage est souvent mesurée par

$$E(\vec{w}) = \frac{1}{2} \sum_{d \in D} \sum_{k \in \text{sorties}} (t_{k,d} - o_{k,d})^2$$

où

- D est l'ensemble de exemples d'apprentissage
- $t_{k,d}$ est la vraie classification de l'instance d pour la sortie k
- $o_{k,d}$ est la valeur de la sortie k pour l'instance d

Backpropagation algorithm

- x_{ij} est l'entrée du noeud j venant du noeud i et $w_{i,j}$ est le poids correspondant.
- δ_n est l'erreur associée à l'unité n . Cette erreur joue le rôle du terme $t - o$ dans le cas d'une seule unité.

```
1 initialise chaque  $w_i$  avec une valeur au hasard
2 tant que l'erreur est trop grande
3   Pour chaque exemple  $(\vec{x}, t) \in D$ 
4     1- calcul de l'état du réseau par propagation
5     2- rétro propagation des erreurs
6       a- pour chaque unité de sortie  $k$ , calcul du terme d'erreur
7          $\delta_k \leftarrow o_k(1 - o_k)(t_k - o_k)$ 
8       b- pour chaque unité cachée  $h$ , calcul du terme d'erreur
9          $\delta_h \leftarrow o_h(1 - o_h) \sum_{k \in \text{sorties}} w_{k,h} \delta_k$ 
10      c- mise à jour des poids  $w_{ji}$ 
11         $w_{j,i} \leftarrow w_{j,i} + \eta \delta_j x_{ji}$ 
```

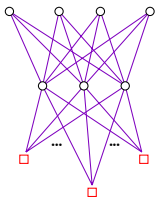
Propriétés

- backpropagation converge vers un minimum local (aucune garantie que le minimum soit global)
- en pratique, il donne de très bon résultat
dans la pratique, le grand nombre de dimensions peut donner des opportunités pour ne pas être bloqué dans un minimum local
- pour le moment, on n'a pas assez de théorie qui explique les bons résultats en pratique
 - ajouter un terme de "momentum"
 - entraîner plusieurs RNA avec les mêmes données mais différents poids de départ (boosting)

Quelles fonctions peut on représenter avec un RNA ?

- fonctions booléennes
- fonctions continues (théoriquement : toute fonction continue peut être approximée avec une erreur arbitraire par un réseau avec deux niveaux d'unités)
- fonctions arbitraires (théoriquement : toute fonction peut être approximée avec une précision arbitraire par un réseau avec 3 niveaux d'unités)

Exemple : reconnaissance de forme



20 personnes, environ 32 photos par personnes, leur tête peut être de face, tournée à gauche, droite, regardant en haut, différentes expressions (content, triste, en colère, neutre), avec ou sans lunettes de soleil. image de 120x128 pixels, noire et blanc, intensité entre 0 et 255.

taux de réussite de 90% pour apprendre l'orientation de la tête et reconnaître une des 20 personnes

