

Cours de Data Science : Réseaux de neurones

Marc Tommasi

24 novembre 2020

Outline

- 1 Petit historique
- 2 Le perceptron
- 3 Perceptron multi couches

Outline

- 1 Petit historique
- 2 Le perceptron
- 3 Perceptron multi couches

Cybernétique : années 50

- Bio-inspiration
- McCulloch-Pitts 43 : modèle du neurone
- Rosenblatt 58/62 : apprendre les poids
- Widrow-Hoff 60 : Adaline comme une instance de la descente de gradient stochastique
- Minsky-Papert 69 : limitation par l'exemple du XOR

Connexionisme : années 80

- une tâche complexe est le résultat de la composition de tâches simples.
- retropropagation du gradient : Rumelhart 86 / Lecun 87
- utilisation de données non étiquetées pour faciliter l'apprentissage de réseaux profonds
- plus trop de connexion avec les neurosciences (on sait peu de choses sur le cerveau, pour l'instant, s'en rapprocher n'a rien donné).

Tirée de Goodfellow

The choice of the functions used to compute these representations is also loosely guided by neuroscientific observations about the functions that biological neurons compute. Modern neural network research, however, is guided by many mathematical and engineering disciplines, and the goal of neural networks is not to perfectly model the brain. It is best to think of feedforward networks as function approximation machines that are designed to achieve statistical generalization, occasionally drawing some insights from what we know about the brain, rather than as models of brain function

Deep learning : depuis milieu des années 2000

- Beaucoup de données (étiquetées)
- Des machines puissantes (GPU, ...)
- Explosion de la taille des réseaux (ver de terre il y a 15 ans, aujourd'hui comparable à la grenouille)
- Résolument inspiré par les statistiques, les probabilités, l'optimisation, le calcul numérique, ...

Évolutions I

Évolution des type de modèles

- Modèle linéaire, un seul neurone (perceptron)
- Modèles à plusieurs couches,
- Réseaux de convolution,
- Modèles non linéaires,
- Modèles récurrents.

Évolution des jeux de données

- de quelques centaines ou milliers d'exemples
- à quelques millions et milliards aujourd'hui
- 10 millions d'exemples étiquetés pour avoir des performances proches de l'homme dans bien des cas. . .

Évolutions II

Évolution des performances

- Des réseaux aujourd'hui avec des millions de neurones
- Impact impressionnant
 - ▶ en vision
 - ▶ en reconnaissance de la parole

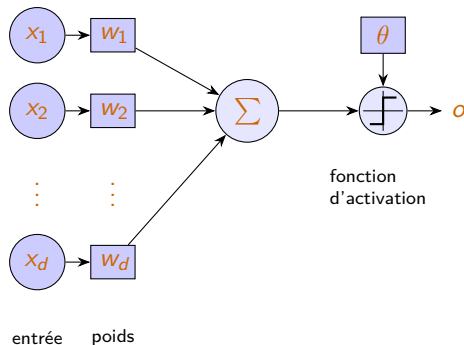
Outline

- 1 Petit historique
- 2 **Le perceptron**
- 3 Perceptron multi couches

Un problème de classification supervisée binaire

- Les données sont décrites par d valeurs réelles : $\mathbf{x} \in \mathbb{R}^d$
- Les résultats sont binaires : on doit associer une valeur binaire y à chaque donnée (y vaut -1 ou 1)
- On dispose de données étiquetées : $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$
- On cherche la meilleure fonction qui permet d'obtenir les résultats (étiquettes y) en fonction des entrées (données \mathbf{x}).

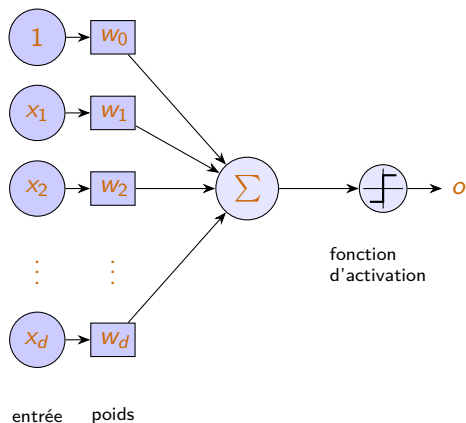
Perceptron linéaire à seuil



- Entrées : d valeurs x_1, \dots, x_d ;
Sortie : o une valeur binaire
- Fonctionnement : d coefficients (synaptiques) w_1, \dots, w_d , le biais θ ,

$$o = \begin{cases} 1 & \text{si } \mathbf{w}^\top \mathbf{x} > \theta \\ -1 & \text{sinon} \end{cases}$$

Perceptron linéaire à seuil



- Entrées : d valeurs x_1, \dots, x_d ;
Sortie : o une valeur binaire
- Fonctionnement : $d + 1$ coefficients (synaptiques)
 w_0, \dots, w_d , entrée $x_0 = 1$,
$$o = \begin{cases} 1 & \text{si } \mathbf{w}^\top \mathbf{x} > 0 \\ -1 & \text{sinon} \end{cases}$$

Exemple (Exercices)

- Quel est le perceptron linéaire à seuil qui calcule un OU logique entre n variables binaires (prenant les valeurs 0 ou 1)
- Donnez une représentation graphique de ce que vous avez proposé dans le cas de deux variables : représentez par des points de deux couleurs les valeurs possibles et la frontière de décision pour la sortie (considérant cette fois que les valeurs des x_i peuvent être réelles).
- Représentez maintenant les points correspondants à une fonction XOR.
- Démontrez qu'un perceptron linéaire à seuil ne peut séparer ces points correspondant au XOR.

Algorithme par correction d'erreur I

Intuition

- considérons des valeurs pour \mathbf{w} , comment les adapter en fonction d'un exemple (\mathbf{x}, y) ?
- rappelons que les poids w_i associés aux entrées x_i qui sont à 0 ne participent pas à la décision.
- Supposons $y = 1$.
 - ▶ Si $\mathbf{w}^\top \mathbf{x} > 0$ alors on n'a rien à faire
 - ▶ Si $\mathbf{w}^\top \mathbf{x} \leq 0$ alors il faut augmenter les poids associés aux entrées x_i qui sont à 1.
- Supposons $y = -1$.
 - ▶ Si $\mathbf{w}^\top \mathbf{x} \leq 0$ alors on n'a rien à faire
 - ▶ Si $\mathbf{w}^\top \mathbf{x} > 0$ alors il faut baisser les poids associés aux entrées x_i qui sont à 1.

Algorithme par correction d'erreur II

L'algorithme d'apprentissage

- Entrée : Un échantillon $\{(\mathbf{x}_1, y_1) \dots, (\mathbf{x}_m, y_m)\}$
- $t \leftarrow 0$
- Initialiser les poids \mathbf{w}_0 à des valeurs aléatoires
- Répéter
 - ▶ Choisir/recevoir un exemple k , (\mathbf{x}_k, y_k) dans S .
 - ▶ Calculer la sortie $\hat{y} = \text{sign}(\mathbf{w}_t^\top \mathbf{x}_k)$
 - ▶ $t \leftarrow t + 1$
 - ▶ Si $(\hat{y} \neq y_k)$ alors $\mathbf{w}_t \leftarrow \mathbf{w}_{t-1} + y_k \mathbf{x}_k$
- Jusqu'à ...
- Sortie : \mathbf{w}_t

Algorithme par correction d'erreur III

Exemple avec le OU

- On commence avec $w_0 = 0$; $w_1 = 1$ et $w_2 = -1$.
- On passe tous les exemples de l'échantillon :
 $S = \{((1, 0, 0), 0), ((1, 0, 1), 1), ((1, 1, 0), 1), ((1, 1, 1), 1)\}.$
- Calculer les sorties successives et les mises à jour

Propriétés de cet algorithme I

Si l'échantillon est linéairement séparable, si les exemples sont présentés équitablement, l'algorithme s'arrête et calcule un séparateur linéaire pour S .

Theorem (Terminaison)

Soit $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathbb{R}^d$ tels que il existe $r > 0$ avec $\|\mathbf{x}_k\| \leq r$ pour tout $k \in [1, m]$.
Si il existe $\rho > 0$ et $\mathbf{v} \in \mathbb{R}^d$, $\|\mathbf{v}\| = 1$, tels que pour tout $k \in [1, m]$,
 $\rho \leq y_k(\mathbf{v}^\top \mathbf{x}_k)$, alors le nombre d'itérations de l'algorithme du perceptron avec $\mathbf{x}_1, \dots, \mathbf{x}_m$ en entrée est borné par r^2/ρ^2 .

- r : est la norme maximale des données
- \mathbf{v} : est l'hyperplan séparateur des données (unitaire)
- ρ : est la *marge*

Propriétés de cet algorithme II

Sketch of proof, borne inférieure

Supposons qu'à l'étape n , l'algorithme a déjà effectué M mises à jour aux temps t_1, t_2, \dots, t_M .

$$\begin{aligned}\|\mathbf{w}_n\|^2 &= \|\mathbf{w}_{t_M} + y_{t_M} \mathbf{x}_{t_M}\|^2 \\ &= \|\mathbf{w}_{t_M}\|^2 + \|\mathbf{x}_{t_M}\|^2 + 2y_{t_M} \mathbf{w}_{t_M}^\top \mathbf{x}_{t_M}\end{aligned}$$

Comme par l'algorithme, $y_{t_M} \mathbf{w}_{t_M}^\top \mathbf{x}_{t_M} < 0$ on a

$$\|\mathbf{w}_n\|^2 \leq \|\mathbf{w}_{t_M}\|^2 + \|\mathbf{x}_{t_M}\|^2$$

.

Par induction

$$\|\mathbf{w}_n\|^2 \leq \|\mathbf{w}_0\|^2 + \sum_{i=1}^M \|\mathbf{x}_{t_i}\|^2 \leq r^2 M$$

Propriétés de cet algorithme III

Sketch of proof, borne supérieure

On a \mathbf{v} qui est unitaire (de norme 1) donc, pour tout vecteur \mathbf{w} , la norme de \mathbf{w} est plus grande que le produit scalaire $\mathbf{v}^\top \mathbf{w}$. Ici,

$$\begin{aligned}\|\mathbf{w}_n\| &\geq \mathbf{v}^\top \mathbf{w}_n \\ &\geq \mathbf{v}^\top (\mathbf{w}_{t_M} + y_{t_M} \mathbf{x}_{t_M}) \\ &\geq \mathbf{v}^\top \mathbf{w}_{t_M} + y_{t_M} \mathbf{v}^\top \mathbf{x}_{t_M} \\ &\geq \mathbf{v}^\top \mathbf{w}_{t_M} + y_{t_M} \rho\end{aligned}$$

Par induction (on suppose qu'on initialise \mathbf{w}_0 à 0)

$$\|\mathbf{w}_n\| \geq M\rho$$

En conclusion on a $\|\mathbf{w}_n\|^2 \geq M^2 \rho^2$ et $\|\mathbf{w}_n\|^2 \leq r^2 M$. Donc

$$M \leq r^2 / \rho^2$$

Propriétés et limitations

- Plus la marge est petite, plus l'algorithme peut prendre du temps à converger.
- On remarque que le vecteur w est calculé avec des combinaisons des exemples x .
- Les poids sont entiers
- Pas de garantie de convergence dans le cas non séparable
- De nombreux efforts pour réaliser des approximations et tolérer des erreurs
- On peut voir l'algorithme comme une descente de gradient stochastique. . .

Rappel descente de gradient

Algo générique

- Entrée : Un échantillon $S = \{(\mathbf{x}_1, y_1) \dots, (\mathbf{x}_m, y_m)\}$, une fonction de perte f
 - Initialiser les poids \mathbf{w}_0 à des valeurs aléatoires
 - Pour t de 1 à T
 - ▶ Prendre \mathbf{x}, y dans S
 - ▶ $\mathbf{w}_t \leftarrow \mathbf{w}_{t-1} - \eta \nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x}, y)$
 - Sortie : \mathbf{w}_T
-
- De nombreuses fonctions de perte peuvent être utilisées selon les cas (MSE (régression linéaire etc. . .), logistique (régression logistique), exponentielle (dans Adaboost), Hinge loss
 - Hinge Loss : $\max(0, 1 - y_k \mathbf{w}^\top \mathbf{x}_k)$
 - ▶ Dérivée : $-y_k \mathbf{x}_k$ si les signes de y_k et $\mathbf{w}^\top \mathbf{x}_k$ diffèrent
 - ▶ On retrouve l'algorithme par correction d'erreur avec $\eta = 1$

Outline

- 1 Petit historique
- 2 Le perceptron
- 3 **Perceptron multi couches**

Description

- La sortie des neurones d'une couche sont l'entrée des neurones de la couche suivante
- On ne crée pas de cycle
- Le réseau réalise des compositions de plusieurs fonctions
- Ces réseaux sont appelés feed-forward

Petites question

- Peut on faire une porte NAND avec un perceptron ?
- Que peut-on faire si on compose des portes NAND ?

Autres fonctions d'activation

- La fonction d'activation à seuil fait un changement très brutal dans un réseau
- Ce changement rend l'apprentissage très difficile : de petits changements ont de gros effets !
- Passage à des fonctions plus « lisses »
 - ▶ la sigmoïde par exemple.

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

la dérivée

$$\sigma'(x) = \sigma(x)(1 - \sigma(x))$$

- ▶ La tangente hyperbolique (tanh) ou la hinge/ReLU.

