

## PROJET L3

But : une petite introduction au Deep Learning avec l'exemple d'un perceptron multicouche.

### 1. INTRODUCTION

Nous allons programmer avec numpy un petit perceptron multicouche dont le but est de faire de la classification à partir d'un jeu de données très simple. Voici un exemple de tel jeu de données (Fig. 1).

On dispose de dix points dans le carré  $[0, 1] \times [0, 1]$ . Ce sont les données qui vont servir pour l'apprentissage. Elles sont classées en deux classes : les bons sont les ronds verts, les méchants les triangles bleus (c'est comme ça la vie, il y a les bons et les méchants...). Le type de réseau que l'on considère ici, le perceptron multicouche, nécessite d'avoir de telles données annotées.

Le but du réseau est de nous fournir, à partir de cet ensemble de données annotées, une frontière de décision entre le bien et le mal... Une fois connue cette frontière, nous pourrions classer les nouvelles données apparaissant en bonnes ou méchantes... Voici la frontière que l'on obtient sur cet exemple (Fig. 1).

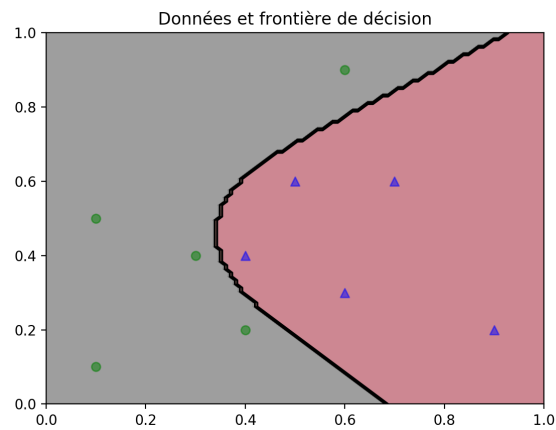


FIGURE 1. Les données et la frontière de décision

## 2. LES MATHÉMATIQUES

L'idée générale du perceptron est de considérer des neurones fabriqués à partir d'équations d'hyperplans et d'une fonction d'activation. Les hyperplans sont donnés par :

$$(1) \quad Wx + b = 0$$

où  $x$  est un vecteur de  $\mathbb{R}^2$ ,  $W$  est une matrice  $2 \times 2$  et  $b$  est un vecteur de  $\mathbb{R}^2$ . Dans la terminologie  $W$  contient les poids, et  $b$  contient les biais.

L'activation du neurone se fait suivant la valeur de

$$(2) \quad \sigma(Wx + b)$$

où  $\sigma$  est la fonction d'activation considérée. Dans notre cas, nous utiliserons la sigmoïde

$$(3) \quad \sigma(x) = \frac{1}{1 + e^{-x}}$$

L'un des avantages de la sigmoïde est la forme très simple de l'équation différentielle qu'elle

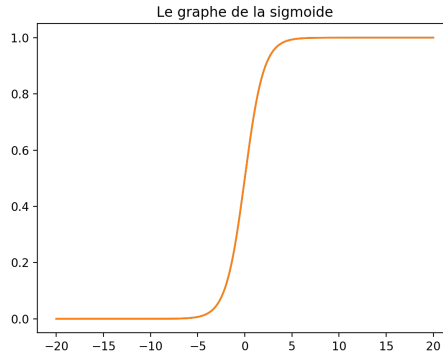


FIGURE 2. Le graphe de la sigmoïde

satisfait. On vérifie facilement que

$$(4) \quad \sigma'(x) = \sigma(x)(1 - \sigma(x))$$

Dans l'équation (2), les poids,  $W$ , et les biais,  $b$ , sont des paramètres que le réseau va se charger de trouver par optimisation.

Voici le petit, très petit, réseau que l'on va programmer. Il contient une couche d'entrée, une couche de sortie, et deux couches dites cachées.

Les données d'entrée sont des vecteurs  $x$  de  $\mathbb{R}^2$ . La sortie du layer 2 est donnée par

$$(5) \quad \sigma(W^{[2]}x + b^{[2]}) \in \mathbb{R}^2$$

où  $W^{[2]}$  est une matrice  $2 \times 2$  et  $b^{[2]}$  est un vecteur de  $\mathbb{R}^2$ . La sortie est donc un vecteur de  $\mathbb{R}^2$ .

Le layer 3 contient trois neurones. Sa sortie sera un vecteur de  $\mathbb{R}^3$  donnée par

$$(6) \quad \sigma\left(W^{[3]}\sigma(W^{[2]}x + b^{[2]}) + b^{[3]}\right) \in \mathbb{R}^3$$

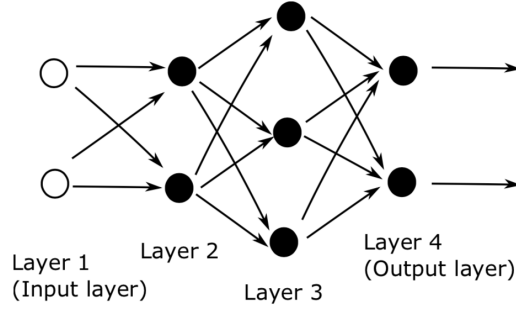


FIGURE 3. Le perceptron multicouche

où  $W^{[3]}$  est une matrice  $3 \times 2$  et  $b^{[3]}$  est un vecteur de  $\mathbb{R}^3$ .

Pour le layer 4, on a pour sortie le vecteur de  $\mathbb{R}^2$

$$(7) \quad \sigma \left( W^{[4]} \sigma \left( W^{[3]} \sigma(W^{[2]}x + b^{[2]}) + b^{[3]} \right) + b^{[4]} \right) \in \mathbb{R}^2$$

où  $W^{[4]}$  est une matrice  $2 \times 3$  et  $b^{[4]}$  est un vecteur de  $\mathbb{R}^2$ .

On considère alors la fonction  $F : \mathbb{R}^2 \longrightarrow \mathbb{R}^2$  définie par

$$(8) \quad F(x) = \sigma \left( W^{[4]} \sigma \left( W^{[3]} \sigma(W^{[2]}x + b^{[2]}) + b^{[3]} \right) + b^{[4]} \right)$$

Cette fonction dépend des paramètres suivants : les poids

$$w_{ij}^{[2]}, w_{ij}^{[3]}, w_{ij}^{[4]}$$

et les biais

$$b_i^{[2]}, b_i^{[3]}, b_i^{[4]}$$

La dimension de l'espace des paramètres est donc 23 (c'est très petit comparé aux usines à gaz actuelles pour lesquelles cette dimension peut être de l'ordre du million...).

Cette fonction est liée à notre problème de classification de la façon suivante. Notons  $x^i$ ,  $i = 1, 2, \dots, 10$  les données. On pose

$$(9) \quad y(x^i) = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

si  $x^i$  est bon et

$$(10) \quad y(x^i) = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$

si  $x^i$  est méchant. La fonction de coût que l'on considère (cost function ou encore loss) est donnée par

$$(11) \quad \text{Cost} \left( W^{[2]}, W^{[3]}, W^{[4]}, b^{[2]}, b^{[3]}, b^{[4]} \right) = \frac{1}{10} \sum_{i=1}^{10} \frac{1}{2} \|y(x^i) - F(x^i)\|_2^2$$

Le but est de minimiser cette fonction Cost qui est une fonction des paramètres du réseau. On trouve ainsi des poids et des biais que l'on ré-injecte dans  $F$ . La frontière de décision est alors donnée par  $F_1(x) > F_2(x)$  avec  $F(x) = (F_1(x), F_2(x))$ .

### 3. LA DESCENTE DE GRADIENT

Ce type de problème de minimisation fait appel à une descente de gradient : c'est le coeur du Deep Learning. Les fonctions coût à minimiser peuvent avoir un très grand nombre de minima locaux. L'étude de la convergence des descentes est un problème particulièrement étudié à l'heure actuelle. Les Figures 3 et 4 montrent une configuration de données pour lesquelles la descente de gradient a stagné pendant 3 millions d'itérations avant de décider à se remettre à descendre...

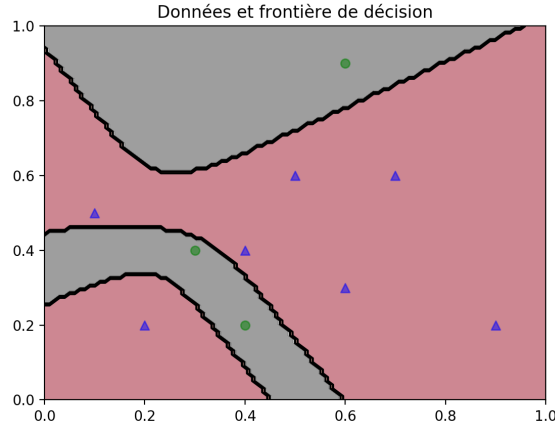


FIGURE 4. Une configuration particulière

On a

$$(12) \quad \text{Cost}(p + \Delta p) \simeq \text{Cost}(p) + \sum_{r=1}^d \frac{\partial \text{Cost}(p)}{\partial p_r} \Delta p_r$$

avec  $\text{Cost} : \mathbb{R}^d \rightarrow \mathbb{R}$ ,  $d$  est la dimension des paramètres,  $p = (p_1, p_2, \dots, p_d)$ . Par définition

$$(13) \quad \langle \nabla \text{Cost}(p), v \rangle = d_p \text{Cost}(v)$$

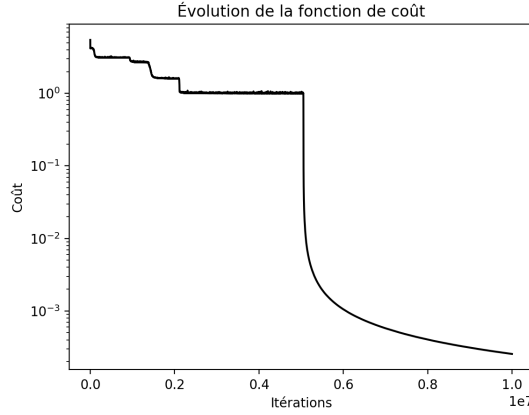


FIGURE 5. L'évolution de la fonction coût correspondante

pour tout vecteur  $v$  de  $\mathbb{R}^d$ . On peut considérer dans notre contexte (euclidien) que  $\nabla \text{Cost}(p)$  est le vecteur de  $\mathbb{R}^d$  de composantes

$$(14) \quad (\nabla \text{Cost}(p))_r = \frac{\partial \text{Cost}(p)}{\partial p_r}$$

L'équation (12) peut donc s'écrire

$$(15) \quad \text{Cost}(p + \Delta p) \simeq \text{Cost}(p) + \nabla \text{Cost}(p)^T \Delta p$$

Pour faire diminuer le plus rapidement la fonction  $\text{Cost}$ , on choisit  $\Delta p$  de sorte que

$$\nabla \text{Cost}(p)^T \Delta p$$

soit le plus négatif possible. D'après Cauchy - Schwarz, il faut choisir

$$\Delta p = -\eta \nabla \text{Cost}(p)$$

où  $\eta$  est un pas appelé le Learning Rate. La mise à jour est donnée par

$$(16) \quad p \longleftarrow p - \eta \nabla \text{Cost}(p)$$

Dans notre cas, comme dans la plupart des cas, notre fonction coût est une somme. On a donc

$$(17) \quad \nabla \text{Cost}(p) = \frac{1}{N} \sum_{i=1}^N \nabla \text{Cost}_i(p)$$

avec  $N$  le nombre de données (ici 10) et

$$(18) \quad \text{Cost}_i = \frac{1}{2} \|y(x^i) - F(x^i)\|_2^2$$

En général, le nombre  $N$  est très grand ce qui pose des problèmes pour le calcul du gradient complet. Pour pallier cette difficulté, on utilise un gradient dit stochastique en expliquant comment sélectionner un nombre restreint de fonctions  $\text{Cost}_i$  et ce de façon à

couvrir au mieux les données. Ici, par soucis de simplicité, on va sélectionner une seule fonction coût  $\text{Cost}_i$  en tirant, pour chaque mise à jour, un indice  $i$  de façon aléatoire et uniforme dans l'ensemble  $\{1, 2, \dots, N = 10\}$  :

$$p \leftarrow p - \eta \nabla \text{Cost}_i(p)$$

#### 4. LA BACKPROP

On commence par introduire quelques notations plus génériques pour les calculs. On note  $x \in \mathbb{R}^{n_1}$  les données d'entrée (pour nous  $n_1 = 2$ ). On note également  $a_j^{[l]}$  la sortie du  $j$ -ème neurone du layer  $l$ . On a donc

$$(19) \quad a^{[1]} = x \in \mathbb{R}^{n_1}$$

$$(20) \quad a^{[l]} = \sigma \left( W^{[l]} a^{[l-1]} + b^{[l]} \right) \in \mathbb{R}^{n_l}$$

pour  $l = 2, 3, \dots, L$  (pour nous  $L = 4$  et  $n_4 = 2$ ). On aura aussi besoin de l'entrée du  $j$ -ème neurone du layer  $l$  (avant activation)

$$(21) \quad z_j^{[l]} = \left( W^{[l]} a^{[l-1]} + b^{[l]} \right)_j$$

Notre problème se résume donc à calculer le gradient de la fonction

$$(22) \quad C = \frac{1}{2} \|y - a^{[L]}\|_2^2$$

vue comme fonction des poids et des biais (qui interviennent dans la sortie  $a^{[L]}$  du réseau).

La propagation de l'information dans le réseau est donnée par

$$(23) \quad a^{[l]} = \sigma \left( z^{[l]} \right)$$

pour  $l = 2, 3, \dots, L$ .

On note  $\delta^{[l]} \in \mathbb{R}^{n_l}$  le vecteur de composantes

$$(24) \quad \delta_j^{[l]} = \frac{\partial C}{\partial z_j^{[l]}}$$

pour  $j = 1, 2, \dots, n_l$  et  $l = 2, 3, \dots, L$ .

Pour faire les calculs on peut utiliser le produit de Hadamard : si  $x$  et  $y$  sont deux vecteurs de  $\mathbb{R}^n$ , le produit de Hadamard  $x \circ y$  de  $x$  et  $y$  est le vecteur de  $\mathbb{R}^n$  de composantes  $(x \circ y)_i = x_i y_i$ , pour tout  $i = 1, 2, \dots, n$ . Tout est dans le résultat suivant.

**Proposition -** On a

$$(25) \quad \delta^{[L]} = \sigma' \left( z^{[L]} \right) \circ \left( a^{[L]} - y \right)$$

$$(26) \quad \delta^{[l]} = \sigma' \left( z^{[l]} \right) \circ \left( W^{[l+1]} \right)^T \delta^{[l+1]}$$

pour  $2 \leq l \leq L - 1$  ;

$$(27) \quad \frac{\partial C}{\partial b_j^{[l]}} = \delta_j^{[l]}$$

pour  $2 \leq l \leq L$  ;

$$(28) \quad \frac{\partial C}{\partial w_{jk}^{[l]}} = \delta_j^{[l]} a_k^{[l-1]}$$

pour  $2 \leq l \leq L$ , en notant  $w_{jk}^{[l]}$  les coefficients de la matrice  $W^{[l]}$ .

**Preuve - Exercice !!**

On peut encore condenser un peu l'écriture en introduisant la matrice  $D^{[l]}$  de taille  $n_l \times n_l$  diagonale dont les termes diagonaux sont

$$D_{ii}^{[l]} = \sigma' \left( z_i^{[l]} \right)$$

On a alors

$$\delta^{[L]} = D^{[L]} \left( a^{[L]} - y \right)$$

et

$$\delta^{[l]} = D^{[l]} \left( W^{[l+1]} \right)^T \delta^{[l+1]}$$

L'écriture finale est donc

$$(29) \quad \delta^{[l]} = D^{[l]} \left( W^{[l+1]} \right)^T D^{[l+1]} \left( W^{[l+2]} \right)^T \dots D^{[L-1]} \left( W^{[L]} \right)^T D^{[L]} \left( a^{[L]} - y \right)$$

Il ne faut pas oublier que le calcul de la dérivée de la sigmoïde est trivial. L'équation (29) est l'équation de la *Back-propagation*.

Le perceptron multicouche est l'exemple type d'un réseau *Forward/Backward*. La passe Forward permet d'évaluer les fonctions d'activation, les poids et la sortie  $a^{[L]}$ . La passe Backward permet de calculer l'erreur et la mise à jour. Ce double mécanisme doit bien apparaître dans le programme.