

# Réseaux de neurones (Notes de cours)

Romain Tavenard



# Chapitre 1

## Préambule

Ce document est un ensemble de notes associées au module de classification supervisée pour la deuxième année de master de Statistiques de l'Université de Rennes 2, et plus particulièrement la partie sur les réseaux de neurones. Il est distribué librement (sous licence [CC BY-NC-SA](#) plus précisément) et se veut évolutif, n'hésitez donc pas à faire vos remarques à son auteur dont vous trouverez le contact sur [sa page web](#).

### 1.1 Contenu du cours



## Chapitre 2

# Régression et perceptron

### 2.1 Retour sur la régression linéaire

Dans le cas d'une régression aux moindres carrés ordinaires, le modèle est le suivant :

$$y_i = \beta_0 + \sum_j \beta_j x_{i,j} + \epsilon_i,$$

où les  $\epsilon_i$  sont i.i.d. gaussiens de moyenne nulle.

On note

$$\hat{y}_i = \beta_0 + \sum_j \beta_j x_{i,j}$$

et l'on souhaite estimer les paramètres  $\beta_j$  minimisant la quantité  $L(\beta) = \frac{1}{2} \sum_i (y_i - \hat{y}_i)^2$ .

Vous connaissez probablement une forme explicite permettant d'exprimer les  $\beta_j$  en fonction des  $y_i$  et des  $x_{i,j}$ . Si vous ne connaissiez pas cette forme explicite (ce sera notre cas pour à peu près tous les réseaux de neurones que nous verrons dans la suite), un moyen d'estimer les  $\beta_j$  serait de minimiser  $L(\beta)$  par descente de gradient. Pour cela, il nous faudra calculer les dérivées de  $L$  par rapport aux différents  $\beta_j$  :

$$\frac{\partial L}{\partial \beta_0} = - \sum_i (y_i - \hat{y}_i) \tag{2.1}$$

$$\forall j \geq 1, \frac{\partial L}{\partial \beta_j} = - \sum_i (y_i - \hat{y}_i) x_{i,j} \tag{2.2}$$

**Travail personnel:** Reprendre les formules précédentes pour le cas d'une régression logistique binaire ( $y_i$  peut prendre les valeurs 1 ou -1).

Dans la suite, on aura parfois recours à cet exemple de la régression aux moindres carrés ordinaires pour comprendre comment sont estimés les paramètres des réseaux de neurones. L'extension à des problèmes de classification (binaire ou non) revient alors à effectuer le même type de modifications que celles présentes dans ce travail personnel (changer l'expression de la vraisemblance, passer à la log-vraisemblance, recalculer les dérivées).

## 2.2 Le perceptron : notations et représentation

Avant de présenter ce qu'est un réseau de neurones, nous allons nous intéresser à un modèle dans lequel on n'a qu'un neurone : le perceptron. Dans ce modèle, on suppose :

$$y_i = f \left( \beta_0 + \sum_j \beta_j x_{i,j} \right) + \epsilon_i$$

et on note

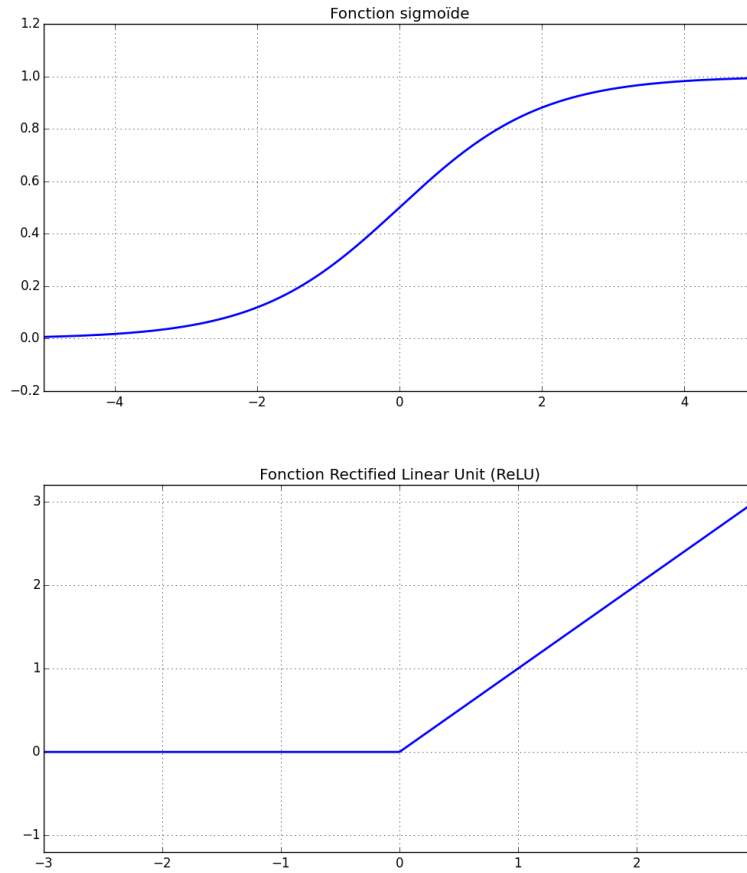
$$\hat{y}_i = f \left( \beta_0 + \sum_j \beta_j x_{i,j} \right).$$

Les paramètres de ce modèle sont les  $\beta_j$  et la fonction  $f$  est appelée fonction d'activation.

Pour cette fonction d'activation, plusieurs choix sont possibles. On peut citer pour exemples les fonctions sigmoïde et *ReLU* (Rectified Linear Unit) :

$$f_1(x) = \frac{1}{1 + e^{-x}}$$

$$f_2(x) = ReLU(x) = \begin{cases} x, & \text{si } x \geq 0 \\ 0, & \text{sinon} \end{cases}$$



Puisque nous en aurons besoin par la suite, nous pouvons d'ores et déjà calculer la dérivée de ces fonctions :

$$\frac{\partial f_1}{\partial x} = f_1(x) \cdot (1 - f_1(x))$$

$$\frac{\partial f_2}{\partial x} = \begin{cases} 1, & \text{si } x > 0 \\ 0, & \text{sinon} \end{cases}$$

Notez que  $f_2$  n'est pas dérivable en 0. On choisit par convention de prolonger sa dérivée en 0 par la valeur 0. Il faut comprendre ici que ce calcul de dérivée nous servira à effectuer notre descente de gradient. En pratique, la probabilité que l'on ait à évaluer  $\frac{\partial f_2}{\partial x}$  en 0 est nulle et donc ce cas ne se présentera pas à nous.

Considérons une fois encore le problème de régression aux moindres carrés ordinaires avec ce modèle. On cherche donc à minimiser  $L(\beta) = \frac{1}{2} \sum_i (y_i - \hat{y}_i)^2$  et l'on n'a, cette fois, plus d'expression analytique pour les  $\beta_j$ . On va donc chercher à effectuer une descente de gradient.

Pour cela, on doit calculer  $\frac{\partial L}{\partial \beta_j}$  pour  $j = 0 \dots D$  où  $D$  est la dimension de l'espace dans lequel vivent nos  $x_i$ . On sépare le cas  $j = 0$  et on obtient :

$$\frac{\partial L}{\partial \beta_0} = 1$$

### 2.2.1 Poids

### 2.2.2 Fonction d'activation

### 2.2.3 Calcul de gradient

Blabla je parle de l'équation (2.3) qui est super.

$$y = mx + b \tag{2.3}$$



## Chapitre 3

# Perceptron multi-couches

<http://cs231n.github.io/neural-networks-1/>

### 3.1 Cas d'un réseau de neurones à une couche cachée

#### 3.1.1 Calcul de gradient

#### 3.1.2 Visualisation de frontières apprises

Exemple intéressant : <http://cs231n.github.io/neural-networks-case-study/>

### 3.2 Cas multi-couches

#### 3.2.1 Calcul de gradient

<http://cs231n.github.io/optimization-2/>

#### 3.2.2 Visualisation

Comprendre qu'on apprend une représentation (eg PCA améliorée) suivie d'un classifieur linéaire (logistic regression)

- <https://rajarsheem.wordpress.com/2017/05/04/neural-networks-dynamics/>
- <http://colah.github.io/posts/2014-03-NN-Manifolds-Topology/>

### 3.2.3 Exemples adversaires

## 3.3 Quelques considérations pratiques

### 3.3.1 Initialisation des poids du réseau

### 3.3.2 Régularisation

<http://cs231n.github.io/neural-networks-2/>

#### 3.3.2.1 Régularisations L1/L2

#### 3.3.2.2 Utilisation du *dropout*

## Chapitre 4

# Réseaux de neurones convolutionnels

### 4.1 Motivation

### 4.2 Mise en oeuvre

<http://cs231n.github.io/convolutional-networks/>

### 4.3 Apprentissage par transfert et *fine-tuning*

<http://cs231n.github.io/transfer-learning/>



## Chapitre 5

# Mise en oeuvre avec keras

<https://keras.io>.