

# RÉSEAUX DE NEURONES

ADAM SALVAIL-BÉRARD

**RÉSUMÉ** Le but de cet article est de donner une introduction simple, mais formelle, du procédé de construction des réseaux de neurones artificiels. En plus de la définition d'un réseau de neurones, sont données les versions scalaire, vectorielle et par lot des formules utilisées dans l'apprentissage des paramètres du réseau de neurones. Les algorithmes de propagation des données et de rétropropagation du gradient de l'erreur d'apprentissage sont abordés. Finalement, dans le contexte d'optimisation des réseaux de neurones, une méthode exacte de calcul du produit de la matrice hessienne et d'un vecteur est montrée.

## 1 Introduction

Les réseaux de neurones sont des constructions abstraites simulant l'activité d'un réseau de neurones biologique simplifié. Ils sont utilisés en apprentissage automatique pour construire un modèle à partir de données existantes dans le but d'effectuer des prédictions sur de nouvelles données soit à l'aide de régression, dans le cas continu, ou de classification, dans le cas discret. Dans [SB12] se trouve une introduction plus complète à l'apprentissage automatique et à son cousin, le forage de données.

L'apprentissage automatique<sup>1</sup> concerne l'ajustement de paramètres d'une « boîte » qui prend en entrée des données et fournit en sortie un résultat. C'est la métaphore classique décrivant une fonction. En statistiques, la boîte est souvent simple et peut utiliser l'analyse mathématique pour caractériser les paramètres optimaux. Par exemple, la régression classique utilise une « boîte » constituée d'une combinaison linéaire de fonctions élémentaires pertinentes et l'ajustement consiste simplement à optimiser les coefficients de la combinaison linéaire. Le contexte moderne du forage de données a motivé le développement de « boîtes » plus complexes, impossibles à optimiser analytiquement. Les réseaux de neurones sont un exemple de tels modèles complexes qui ont fait leurs preuves dans plusieurs domaines comme la détection de fraudes ou la reconnaissance de caractères manuscrits.

Ce texte s'intéresse donc aux réseaux de neurones, à ouvrir la boîte métaphorique qui est associée à cette classe de modèles. Un réseau de neurones

---

J'aimerais remercier Hugo Larochelle et Jean-Pierre Dussault pour leurs judicieux conseils et leurs relectures de ce documents. Merci également à Guillaume Douville et Vincent Morissette-Thomas pour leur relecture et leurs conseils.

<sup>1</sup>De l'anglais *machine learning*.

est un graphe orienté de nœuds, appelés neurones, souvent disposés en couches. Ces nœuds reçoivent une information de leurs prédécesseurs (les neurones de la couche précédente) et combinent cette information selon une certaine pondération. Chaque nœud possède également un seuil d'activation. Autour de ce seuil, la valeur de sortie est beaucoup plus sensible aux changements de la valeur d'entrée.

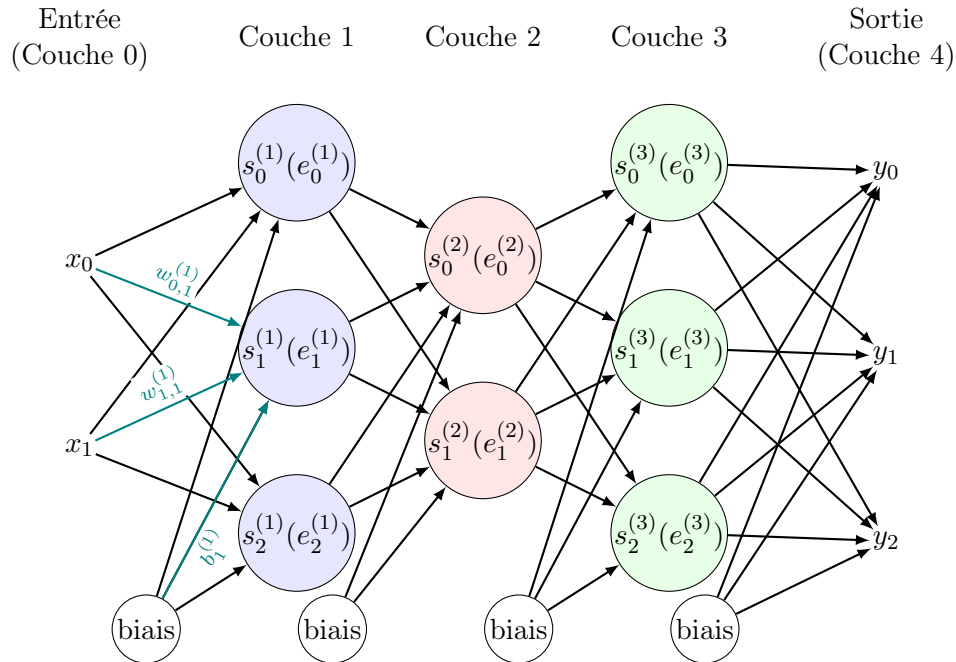


FIGURE 18 – Illustration d'un réseau de neurones *feed-forward* à 5 couches dont 3 couches cachées, une couche d'entrée et une couche de sortie. Ce type de réseau est caractérisé par le fait qu'un neurone ne tire ses entrées que de la couche précédente. À titre d'exemple, les entrées du deuxième nœud de la seconde couche ont été étiquetées en bleu.

Le comportement du réseau de neurones dépend donc de la pondération des apports des nœuds d'une couche à la suivante et des seuils d'activation de chaque nœud. Si ces paramètres sont fixés arbitrairement, les résultats fournis seront médiocres. Il convient d'utiliser un jeu de données nommé ensemble d'entraînement pour optimiser les paramètres afin de minimiser l'erreur de prédiction du réseau. Pour cet ensemble d'entraînement, la prédiction attendue est connue et l'optimisation consiste à minimiser la somme des carrés des différences entre les sorties calculées et les sorties attendues. Contrairement aux modèles classiques en statistique, il est impossible d'obtenir une solution analytique à ces problèmes d'optimisation. Il faut donc optimiser numériquement la fonction d'erreur.

Cet article présente une description rigoureuse de ces réseaux de neurones.

De plus, une méthode d'évaluation des dérivées premières et secondes (dans une direction donnée) d'une fonction d'erreur servant à identifier les paramètres optimaux est proposée. Cette information sur les dérivées est très utile pour les algorithmes d'optimisation qui s'en servent pour trouver rapidement une solution (localement) optimale au problème de minimisation de l'erreur.

## 2 Définitions

Un réseau de neurones agit sur des données. Ces données sont définies formellement comme étant un ensemble de couples  $(\mathbf{x}, \mathbf{c})$  avec  $\mathbf{x} \in \mathbb{R}^n$  et  $\mathbf{c} \in \mathbb{R}^m$ ,  $n, m \in \mathbb{N}^*$ . En outre, un réseau de neurones peut être défini comme étant une fonction  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  qui prend en entrée les informations du vecteur  $\mathbf{x}$  et qui effectue une prédiction  $\mathbf{y} \in \mathbb{R}^m$  du vecteur  $\mathbf{c}$ . Le but de la construction d'un réseau de neurones est de trouver la fonction  $f$  qui donne la meilleure approximation  $\mathbf{y}$  de la cible  $\mathbf{c}$ .

Comme son nom l'indique, un réseau de neurones est une collection de nœuds, nommés neurones, disposés en  $r + 1$  couches. La couche 0 représente l'entrée et possède  $n$  neurones, soit la dimension du vecteur d'entrée  $\mathbf{x}$ . De même, la  $r$ -ième couche, la couche de sortie, possède  $m$  neurones, soit la dimension du vecteur de sortie  $\mathbf{y}$ . Les couches intermédiaires, nommées couches cachées, possèdent un nombre de neurones arbitraires, souvent déterminé par la complexité du problème à résoudre.

Afin de propager les valeurs d'entrée, chaque neurone de la couche  $i$ ,  $0 \leq i \leq r$ , possède une connexion avec tous les neurones de la couche  $i + 1$ <sup>2</sup>. Chacune de ces connexions est pondérée par un poids  $w_{j,k}^{(i)}$  défini comme étant le poids liant la sortie  $s_j^{(i-1)}$  du  $j$ -ième neurone de la couche  $i - 1$  à l'entrée  $e_k^{(i)}$  du  $k$ -ième neurone de la couche  $i$ . Les poids liant la couche  $i - 1$  et la couche  $i$  forment une matrice notée  $\mathbf{w}^{(i)}$ . La Figure 18 illustre ces notations à l'aide d'une représentation graphique d'un réseau de neurones.

Les neurones d'un réseau de neurones représentent les unités de calcul du modèle. Chaque neurone possède une entrée (un scalaire), une fonction d'activation et une sortie. Par exemple, le  $k$ -ième neurone de la couche  $i$  reçoit une valeur d'entrée

$$e_k^{(i)} = \sum_j w_{j,k}^{(i)} s_j^{(i-1)} + b_k^{(i)} \quad (18)$$

et retourne la valeur  $s_k^{(i)} = a_k^{(i)}(e_k^{(i)})$ , soit la fonction d'activation appliquée à l'entrée  $e_k^{(i)}$  du neurone.

La valeur  $b_k^{(i)}$  est un biais ajouté à l'entrée du  $k$ -ième neurone de la couche  $i$  servant à déterminer la sensibilité d'un neurone à son entrée.

<sup>2</sup>Ceci est la représentation classique d'un réseau de neurones. D'autres topologies, c'est-à-dire dispositions des neurones, sont possibles. Dans cet article, seul ce type de réseau de neurones est discuté.

Finalement, il y a la fonction d'erreur

$$E = \frac{1}{2} \|\mathbf{y} - \mathbf{c}\|_2^2 = \frac{1}{2} \sum_{j=1}^m (y_j - c_j)^2, \quad (19)$$

qui dénote l'erreur du réseau de neurones dans son approximation de la sortie  $\mathbf{y}$  par rapport à la cible  $\mathbf{c}$ . C'est cette fonction qui détermine la qualité du modèle.

### 3 Réseau de neurones *feed-forward*

Un réseau de neurones typique est le réseau *feed-forward*. Ce réseau propage l'entrée du réseau vers les couches suivantes sans jamais revenir en arrière. Ce type de réseau est utilisé pour le reste de cet article.

En plus du type de réseau de neurones, il faut également choisir une fonction d'erreur et une fonction d'activation pour les neurones. Ces choix sont souvent guidés par le type de données traitées. Plus de détails sont disponibles dans l'ouvrage de Bishop [Bis06] au chapitre 5.

Dans cet article, la fonction d'activation utilisée est une fonction logistique définie par  $a(x) = \frac{1}{1+\exp(-x)}$ , sauf pour la couche de sortie qui garde la fonction identité comme fonction d'activation. De cette façon, la valeur de sortie

$$s_j^{(i)} = a_j^{(i)}(e_j^{(i)}) = \begin{cases} \frac{1}{1+\exp(-e_j^{(i)})} & \text{pour } i = 0, \dots, r-1 \\ e_j^{(i)} & \text{pour } i = r. \end{cases} \quad (20)$$

Pour la fonction d'erreur, un choix simple est d'utiliser la formule (19), soit la moitié du carré de la distance euclidienne entre la sortie du réseau et la cible.

Avec ces définitions, tous les outils nécessaires à l'établissement d'une prédiction sont présents. Pour effectuer la propagation avec un couple  $(\mathbf{x}, \mathbf{c})$ , les sorties de la couche d'entrée sont assignées aux valeurs du vecteur  $\mathbf{x}$ , c'est-à-dire que  $s_j^{(0)} = x_j$ . Les valeurs sont ensuite propagées dans le réseau de neurones grâce aux formules (18) et (20) pour donner les sorties  $y_j = s_j^{(r)}$  qui sont comparées aux cibles  $c_j$  par la fonction d'erreur  $E$  en (19). Il faut maintenant minimiser la moyenne des erreurs données par la fonction  $E$  sur l'ensemble des données fournies en entrée :

$$E_{\text{moyenne}} = \frac{1}{N} \sum_{t=1}^N E_t$$

où  $N$  est le nombre de couples donnés en entraînement au réseau de neurones et  $E_t$  représente la  $t$ -ième erreur d'apprentissage. Pour y arriver, il suffit d'ajuster les paramètres  $w_{j,k}^{(i)}$  et  $b_j^{(i)}$ .

### 4 Construction du réseau de neurones

Bien entendu, au départ, un réseau de neurones avec des poids initialisés aléatoirement donne de pauvres résultats. C'est pourquoi il y a une étape d'appren-

tissage, c'est-à-dire d'ajustement des poids du réseau afin que ceux-ci puissent donner de meilleures prédictions.

Une façon de voir le problème est d'essayer de minimiser l'erreur de prédiction du modèle. Cette formulation transforme le problème d'apprentissage en un problème d'optimisation pour lequel plusieurs algorithmes existent. L'optimisation de la fonction d'erreur s'effectue par l'ajustement des paramètres de la fonction, c'est-à-dire les poids et les biais du réseau de neurones.

Soit la fonction  $\dim(i)$  qui donne le nombre de neurones de la couche  $i$ . Soit

$$\boldsymbol{\theta} = \begin{bmatrix} w_{j,k}^{(i)} & b_j^{(i)} \end{bmatrix} \in \mathbb{R}^p \quad (21)$$

où  $0 < i \leq r$ ,  $0 \leq j < \dim(i-1)$  et  $0 \leq k < \dim(i)$  qui représentent le vecteur rassemblant chacun des paramètres du réseau de neurones. Le problème d'optimisation associé est de trouver les paramètres du réseau de neurones permettant de minimiser l'erreur  $E(\boldsymbol{\theta})$ . Formellement, il faut trouver

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} \frac{1}{N} \sum_{t=1}^N E_t(\boldsymbol{\theta}). \quad (22)$$

Néanmoins, l'application d'algorithmes d'optimisation nécessite souvent d'avoir plus d'information sur la fonction d'erreur  $E(\boldsymbol{\theta})$ , comme ses dérivées de premier et second ordres par rapport aux paramètres  $\boldsymbol{\theta}$  de la fonction d'erreur.

Trouver les dérivées de la fonction d'erreur est un problème complexe. Puisque la fonction d'erreur, avec toutes les couches du réseau, peut facilement devenir difficile à appréhender, il devient impossible de déduire une forme analytique de la solution au problème de minimisation de l'erreur moyenne. Il faut donc utiliser des méthodes itératives. La méthode la plus élégante pour l'obtention de la dérivée de premier ordre est la rétropropagation de l'erreur de prédiction. Pour la dérivée de second ordre, une méthode proposée par Pearlmutter [Pea94] est décrite. Cette méthode donne le résultat exact du produit de la matrice hessienne et d'un vecteur arbitraire.

#### 4.1 Rétropropagation

L'idée de la rétropropagation est de faire circuler l'information sur la dérivée de la fonction d'erreur à partir de la couche de sortie, où l'erreur de prédiction est connue (la différence entre  $\mathbf{y}$  et  $\mathbf{c}$ ), jusqu'à la couche d'entrée. Pour y arriver, il faut arriver à exprimer la dérivée de la fonction d'erreur d'un nœud en fonction de l'information donnée par les couches suivantes. De cette façon, les valeurs attendues de la couche de sortie pourront être rétropropagées vers l'entrée.

Tout d'abord, la dérivée, à la sortie, de la fonction d'erreur (19) est donnée par

$$\frac{\partial E}{\partial y_i} = y_i - c_i. \quad (23)$$

À partir de cette équation, il est possible de trouver la dérivée de la fonction d'erreur à l'entrée de la couche de sortie. La dérivée reste la même puisque la

fonction d'activation de la dernière couche est la fonction identité, c'est-à-dire

$$\frac{\partial E}{\partial e_j^{(r)}} = \frac{\partial E}{\partial s_j^{(r)}} = \frac{\partial E}{\partial y_j}.$$

À partir de ce moment, il faudra se servir de la formule de dérivée en chaîne pour trouver, de façon récursive, les autres dérivées. En d'autres mots, il faut pouvoir exprimer chaque dérivée de la fonction d'erreur en terme de dérivées venant plus loin dans le réseau de neurones.

Pour ce faire, il faut voir que la fonction d'erreur est affectée par la sortie d'un neurone de la  $i$ -ième couche à travers l'ensemble des neurones de la  $(i + 1)$ -ième couche. Donc,

$$\frac{\partial E}{\partial s_j^{(i)}} = \sum_k w_{j,k}^{(i+1)} \frac{\partial E}{\partial e_k^{(i+1)}}. \quad (24)$$

De plus, la fonction d'erreur est affectée par l'entrée d'un neurone de la  $i$ -ième couche seulement au travers de la sortie de ce même neurone. Il en résulte que

$$\frac{\partial E}{\partial e_j^{(i)}} = \frac{\partial E}{\partial s_j^{(i)}} \frac{\partial s_j^{(i)}}{\partial e_j^{(i)}}.$$

Sachant que la dérivée de la fonction d'activation logistique est donnée par  $\frac{\partial a}{\partial x} = a(1 - a)$ , la dérivée de la fonction d'erreur par l'entrée d'un neurone est donc

$$\frac{\partial E}{\partial e_j^{(i)}} = \frac{\partial E}{\partial s_j^{(i)}} s_j^{(i)} (1 - s_j^{(i)}). \quad (25)$$

La formule (25) permet de rétropropager la dérivée de l'erreur jusqu'à la première couche. À l'aide de ces valeurs, il est possible d'arriver à la dérivée de la fonction d'erreur par rapport à chacun des poids du réseau de neurones, en appliquant la règle de dérivée en chaîne. Donc,

$$\frac{\partial E}{\partial w_{j,k}^{(i)}} = \frac{\partial E}{\partial e_k^{(i)}} \frac{\partial e_k^{(i)}}{\partial w_{j,k}^{(i)}} = \frac{\partial E}{\partial e_k^{(i)}} s_j^{(i-1)}. \quad (26)$$

Finalement, comme les biais sont également des paramètres impliqués dans l'apprentissage d'un réseau de neurones, il faut aussi pouvoir leur trouver une valeur optimale. Il est donc nécessaire de trouver la dérivée de l'erreur par rapport aux biais. De la même façon qu'avec les poids, un biais n'affecte l'erreur de prédiction qu'au travers de l'entrée du nœud auquel il est rattaché. Donc, grâce à la dérivée en chaîne et la formule (18),

$$\frac{\partial E}{\partial b_j^{(i)}} = \frac{\partial E}{\partial e_j^{(i)}} \frac{\partial e_j^{(i)}}{\partial b_j^{(i)}} = \frac{\partial E}{\partial e_j^{(i)}}. \quad (27)$$

## 4.2 Calcul de la dérivée seconde

Pour pouvoir minimiser la fonction d'erreur de manière efficace, plusieurs algorithmes d'optimisation nécessitent des informations sur la dérivée seconde qui représente le comportement de la fonction autour d'un point. Dans le cas de la fonction d'erreur, une fonction avec un domaine dans  $\mathbb{R}^p$ , sa dérivée seconde est en fait la matrice hessienne

$$\mathbf{H} = \begin{bmatrix} \frac{\partial^2 E}{\partial \theta_1^2} & \frac{\partial^2 E}{\partial \theta_1 \partial \theta_2} & \cdots & \frac{\partial^2 E}{\partial \theta_1 \partial \theta_p} \\ \frac{\partial^2 E}{\partial \theta_2 \partial \theta_1} & \frac{\partial^2 E}{\partial \theta_2^2} & \cdots & \frac{\partial^2 E}{\partial \theta_2 \partial \theta_p} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 E}{\partial \theta_p \partial \theta_1} & \frac{\partial^2 E}{\partial \theta_p \partial \theta_2} & \cdots & \frac{\partial^2 E}{\partial \theta_p^2} \end{bmatrix}.$$

Malheureusement, cette matrice, pour de grandes dimensions, devient difficile à stocker et à calculer, même pour un ordinateur. Puisque la plupart des algorithmes utilisent plutôt le produit de la matrice hessienne avec un vecteur arbitraire  $\mathbf{v}$  (résultant en un vecteur de dimension  $p$ ), il est plus facile de calculer directement ce vecteur plutôt que la matrice hessienne entière.

En se basant sur le principe de la propagation de valeurs exposé plus haut, il est possible de trouver la valeur de  $\mathbf{H}\mathbf{v}$  en apportant quelques changements aux formules déjà données. Cette technique, sous la forme présentée dans cet article, est attribuée à Pearlmutter [Pea94].

Le produit hessienne-vecteur peut être défini par

$$\mathbf{H}\mathbf{v} = \left. \frac{\partial}{\partial q} \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta} + q\mathbf{v}) \right|_{q=0}.$$

De cette définition, il est possible d'extraire un nouvel opérateur

$$\mathcal{R}\{f(\boldsymbol{\theta})\} = \left. \frac{\partial}{\partial q} f(\boldsymbol{\theta} + q\mathbf{v}) \right|_{q=0}$$

où  $\mathbf{v}$  représente le vecteur par lequel la matrice hessienne est multipliée et qui est de même dimension que le vecteur  $\boldsymbol{\theta}$ . Pour fins de concordance, les éléments  $v_{j,k}^{(i)}$  et  $v_j^{(i)}$  du vecteur  $\mathbf{v}$  sont respectivement à la même position que les éléments  $w_{j,k}^{(i)}$  et  $b_j^{(i)}$  du vecteur  $\boldsymbol{\theta}$ .

Par définition,  $\mathcal{R}\{w_{j,k}^{(i)}\} = v_{j,k}^{(i)}$  et  $\mathcal{R}\{b_k^{(i)}\} = v_k^{(i)}$  puisque les  $w_{j,k}^{(i)}$  et les  $b_k^{(i)}$  sont les paramètres à optimiser. En appliquant les règles usuelles de dérivée et de dérivée en chaîne, il est possible d'obtenir des formules de propagation pour les calculs de  $\mathcal{R}\{\cdot\}$ . Pour l'entrée dans un neurone, la formule devient, avec la

règle de la dérivée d'un produit,

$$\begin{aligned}
\mathcal{R} \left\{ e_k^{(i)} \right\} &= \mathcal{R} \left\{ \sum_j w_{j,k}^{(i)} s_j^{(i-1)} + b_k^{(i)} \right\} \\
&= \sum_j w_{j,k}^{(i)} \mathcal{R} \left\{ s_j^{(i-1)} \right\} + \mathcal{R} \left\{ w_{j,k}^{(i)} \right\} s_j^{(i-1)} + \mathcal{R} \left\{ b_k^{(i)} \right\} \\
&= \sum_j w_{j,k}^{(i)} \mathcal{R} \left\{ s_j^{(i-1)} \right\} + v_{j,k}^{(i)} s_j^{(i-1)} + v_k^{(i)}.
\end{aligned} \tag{28}$$

Pour la sortie d'un neurone, la formule devient

$$\mathcal{R} \left\{ s_j^{(i)} \right\} = s_j^{(i)} (1 - s_j^{(i)}) \mathcal{R} \left\{ e_j^{(i)} \right\}. \tag{29}$$

En initialisant  $\mathcal{R} \left\{ s_j^{(0)} \right\}$  au vecteur nul et à partir des formules (28) et (29), il est possible de faire la propagation des valeurs données en entrée pour le calcul de la dérivée seconde.

À l'image de la rétropropagation, il faut également rétropropager la dérivée de la fonction d'erreur sur laquelle est appliqué l'opérateur  $\mathcal{R} \{ \cdot \}$ . Il faut donc appliquer l'opérateur  $\mathcal{R} \{ \cdot \}$  à la formule de la dérivée de la sortie :

$$\mathcal{R} \left\{ \frac{\partial E}{\partial y_j} \right\} = \mathcal{R} \left\{ y_j - c_j \right\} = \mathcal{R} \left\{ y_j \right\} = \mathcal{R} \left\{ s_j^{(r)} \right\}. \tag{30}$$

De façon identique à la rétropropagation,

$$\mathcal{R} \left\{ \frac{\partial E}{\partial y_j} \right\} = \mathcal{R} \left\{ \frac{\partial E}{\partial s_j^{(r)}} \right\} = \mathcal{R} \left\{ \frac{\partial E}{\partial e_j^{(r)}} \right\}.$$

Les formules récursives de la rétropropagation deviennent

$$\begin{aligned}
\mathcal{R} \left\{ \frac{\partial E}{\partial e_j^{(i)}} \right\} &= \mathcal{R} \left\{ \frac{\partial E}{\partial s_j^{(i)}} \frac{\partial s_j^{(i)}}{\partial e_j^{(i)}} \right\} \\
&= \frac{\partial E}{\partial s_j^{(i)}} \mathcal{R} \left\{ \frac{\partial s_j^{(i)}}{\partial e_j^{(i)}} \right\} + \mathcal{R} \left\{ \frac{\partial E}{\partial s_j^{(i)}} \right\} \frac{\partial s_j^{(i)}}{\partial e_j^{(i)}} \\
&= \frac{\partial E}{\partial s_j^{(i)}} \frac{\partial^2 s_j^{(i)}}{\partial e_j^{(i)2}} \mathcal{R} \left\{ e_j^{(i)} \right\} + \mathcal{R} \left\{ \frac{\partial E}{\partial s_j^{(i)}} \right\} \frac{\partial s_j^{(i)}}{\partial e_j^{(i)}}
\end{aligned} \tag{31}$$

$$\text{où } \frac{\partial s_j^{(i)}}{\partial e_j^{(i)}} = s_j^{(i)} (1 - s_j^{(i)}) \text{ et } \frac{\partial^2 s_j^{(i)}}{\partial e_j^{(i)2}} = s_j^{(i)} (1 - s_j^{(i)}) (1 - 2s_j^{(i)}).$$



De même, pour la sortie d'un neurone,

$$\begin{aligned}
\mathcal{R} \left\{ \frac{\partial E}{\partial s_j^{(i)}} \right\} &= \mathcal{R} \left\{ \sum_j w_{j,k}^{(i+1)} \frac{\partial E}{\partial e_k^{(i+1)}} \right\} \\
&= \sum_j \left( w_{j,k}^{(i+1)} \mathcal{R} \left\{ \frac{\partial E}{\partial e_k^{(i+1)}} \right\} + \mathcal{R} \left\{ w_{j,k}^{(i+1)} \right\} \frac{\partial E}{\partial e_k^{(i+1)}} \right) \\
&= \sum_j \left( w_{j,k}^{(i+1)} \mathcal{R} \left\{ \frac{\partial E}{\partial e_k^{(i+1)}} \right\} + v_{j,k}^{(i+1)} \frac{\partial E}{\partial e_k^{(i+1)}} \right). \tag{32}
\end{aligned}$$

Finalement, pour les poids

$$\begin{aligned}
\mathcal{R} \left\{ \frac{\partial E}{\partial w_{j,k}^{(i)}} \right\} &= \mathcal{R} \left\{ s_j^{(i-1)} \frac{\partial E}{\partial e_k^{(i)}} \right\} \\
&= s_j^{(i-1)} \mathcal{R} \left\{ \frac{\partial E}{\partial e_k^{(i)}} \right\} + \mathcal{R} \left\{ s_j^{(i-1)} \right\} \frac{\partial E}{\partial e_k^{(i)}} \tag{33}
\end{aligned}$$

et pour les biais

$$\mathcal{R} \left\{ \frac{\partial E}{\partial b_k^{(i)}} \right\} = \mathcal{R} \left\{ \frac{\partial E}{\partial e_k^{(i)}} \right\}. \tag{34}$$

Après tous ces calculs, il est possible de reconstituer le vecteur du produit hessienne-vecteur

$$\begin{aligned}
\mathbf{H}\mathbf{v} &= \mathcal{R} \{ \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}) \} \\
&= \left[ \mathcal{R} \left\{ \frac{\partial E}{\partial w_{0,0}^{(0)}} \right\} \dots \mathcal{R} \left\{ \frac{\partial E}{\partial w_{\dim(r-1)-1, m-1}^{(r)}} \right\} \mathcal{R} \left\{ b_0^{(0)} \right\} \dots \mathcal{R} \left\{ b_{m-1}^{(r)} \right\} \right]^\top. \tag{35}
\end{aligned}$$

## 5 Version vectorielle

Alors que la version scalaire du réseau de neurones est utile à la compréhension, une version vectorielle est plus concise. La plupart des implémentations de cet algorithme d'apprentissage utilisent d'ailleurs la version vectorielle pour utiliser la puissance de calcul des bibliothèques d'algèbre linéaire.

La notation vectorielle se base sur la notation scalaire à partir de laquelle les indices se rattachant à un neurone particulier ont été omis. Par exemple, l'entrée d'un neurone  $\mathbf{e}^{(i)} = [e_0^{(i)} \ e_1^{(i)} \ \dots \ e_{\dim(i)-1}^{(i)}]^\top$ .

L'utilisation de vecteurs apporte quelques changements dont il faut tenir compte. D'abord, l'ensemble des dérivées par rapport à l'erreur est remplacé par le gradient de cette même fonction. Dans cet article, le gradient d'une fonction  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  est défini comme étant

$$\nabla f(x) = \left[ \frac{\partial f}{\partial x_0} \quad \frac{\partial f}{\partial x_1} \quad \cdots \quad \frac{\partial f}{\partial x_{n-1}} \right]^\top.$$

Une exception notable à cette règle (tiré du livre de Rojas [Roj96]) est la fonction d'activation qui est une fonction de  $\mathbb{R}^n$  dans  $\mathbb{R}^n$ , donc qui nécessite de calculer la matrice jacobienne de la fonction, c'est-à-dire que

$$\mathbf{J}_{\mathbf{e}^{(i)}}(\mathbf{a}^{(i)}) = \begin{bmatrix} \frac{\partial a_0^{(i)}}{\partial e_0^{(i)}} & \frac{\partial a_0^{(i)}}{\partial e_1^{(i)}} & \cdots & \frac{\partial a_0^{(i)}}{\partial e_{\dim(i)-1}^{(i)}} \\ \frac{\partial a_1^{(i)}}{\partial e_0^{(i)}} & \frac{\partial a_1^{(i)}}{\partial e_1^{(i)}} & \cdots & \frac{\partial a_1^{(i)}}{\partial e_{\dim(i)-1}^{(i)}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial a_{\dim(i)-1}^{(i)}}{\partial e_0^{(i)}} & \frac{\partial a_{\dim(i)-1}^{(i)}}{\partial e_1^{(i)}} & \cdots & \frac{\partial a_{\dim(i)-1}^{(i)}}{\partial e_{\dim(i)-1}^{(i)}} \end{bmatrix}.$$

Mais comme la dérivée partielle de la fonction d'activation du  $j$ -ième neurone ne dépend que de l'entrée du même neurone, tous les éléments non diagonaux de la matrice  $\mathbf{J}_{\mathbf{e}^{(i)}}(\mathbf{a}^{(i)})$  sont nuls. Donc, dans le cas où la fonction d'activation est la fonction logistique,

$$\mathbf{J}_{\mathbf{s}^{(i)}} = \mathbf{J}_{\mathbf{a}_{\text{logistique}}^{(i)}} = \begin{bmatrix} s_1^{(i)}(1 - s_1^{(i)}) & 0 & \cdots & 0 \\ 0 & s_2^{(i)}(1 - s_2^{(i)}) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & s_n^{(i)}(1 - s_n^{(i)}) \end{bmatrix}. \quad (36)$$

Sachant ceci, il est plus efficace d'implémenter cette matrice comme un vecteur colonne qui est multiplié terme-à-terme à un autre. Par exemple, la formule (39) multiplie un gradient et une matrice jacobienne. Pour implémenter cette multiplication de manière efficace, il est possible de multiplier le gradient et le vecteur contenant les éléments diagonaux de la matrice jacobienne terme à terme. Il en résulte un nouveau vecteur de la même dimension que le gradient, ce qui correspond à la dimension du produit du gradient et de la matrice jacobienne. Cette manipulation évite de mettre en mémoire tous les zéros qui ne sont pas situés sur la diagonale de la matrice jacobienne, ce qui accélère également les calculs.

Pour profiter de ces avantages, l'opérateur  $\text{diag}(\mathbf{A})$  est utilisé pour représenter le vecteur colonne composé des éléments de la diagonale de la matrice  $\mathbf{A}$ . De plus, l'opérateur  $\odot$  est utilisé pour signifier un produit terme à terme de deux matrices de mêmes dimensions.

## 5.1 Propagation

Pour propager les entrées dans le réseau de neurones, il convient de convertir les formules vues précédemment. Dans ce contexte, une entrée du réseau de neurones

est représentée par le vecteur  $\mathbf{x} = \mathbf{s}^{(0)}$ . Ensuite, pour calculer les entrées des neurones des couches suivantes, il faut transformer la formule (18) en

$$\mathbf{e}^{(i)} = \mathbf{w}^{(i)} \mathbf{s}^{(i-1)} + \mathbf{b}^{(i)}.$$

De même pour la sortie d'un neurone, la formule (20) devient

$$\mathbf{s}^{(i)} = \mathbf{a}^{(i)}(\mathbf{e}^{(i)}) = \left[ \frac{1}{1 + \exp(-e_1^{(i)})} \quad \frac{1}{1 + \exp(-e_2^{(i)})} \quad \cdots \quad \frac{1}{1 + \exp(-e_n^{(i)})} \right]^\top.$$

La fonction d'erreur (19) se simplifie pour donner

$$E(\boldsymbol{\theta}) = \frac{1}{2} \|\mathbf{y} - \mathbf{c}\|_2^2 \quad (37)$$

dont le gradient est tout simplement

$$\nabla_{\mathbf{s}^{(r)}} E = \nabla_{\mathbf{y}} E = \mathbf{y} - \mathbf{c}. \quad (38)$$

Ces fonctions permettent de faire l'ensemble des calculs en faisant abstraction des neurones et en ne manipulant que les vecteurs ou matrices se rattachant aux couches.

## 5.2 Rétropropagation

Pour retrouver la version vectorielle de la rétropropagation, il faut refaire les mêmes opérations que pour la version scalaire à partir de la dernière couche. L'erreur par rapport à la sortie est déjà connue (38) et la fonction d'activation de la dernière couche est une fonction linéaire ( $\mathbf{s}^{(r)} = \mathbf{e}^{(r)}$ ). Puisque la dérivée de cette fonction linéaire donne 1, l'entrée de la dernière couche est égale à sa sortie, c'est-à-dire que  $\nabla_{\mathbf{e}^{(r)}} E = \nabla_{\mathbf{s}^{(r)}} E = \mathbf{y} - \mathbf{c}$ .

À partir de là, il est possible de retrouver la formulation des gradients des autres couches de façon récursive. Le gradient de la sortie s'obtient, à l'image de la formule (24), par

$$\nabla_{\mathbf{s}^{(i)}} E = (\mathbf{w}^{(i+1)})^\top \nabla_{\mathbf{e}^{(i+1)}} E.$$

Quant à lui, le gradient de l'entrée d'un neurone suit le principe de la formule (25) :

$$\nabla_{\mathbf{e}^{(i)}} E = \mathbf{J}_{\mathbf{e}^{(i)}}(\mathbf{s}^{(i)}) \nabla_{\mathbf{s}^{(i)}} E = \text{diag} \left( \mathbf{J}_{\mathbf{e}^{(i)}}(\mathbf{s}^{(i)}) \right) \odot \nabla_{\mathbf{s}^{(i)}} E \quad (39)$$

où  $\mathbf{J}_{\mathbf{e}^{(i)}}(\mathbf{s}^{(i)})$  est définie par la formule (36).

Ne reste qu'à utiliser ces deux vecteurs pour pouvoir calculer la valeur des gradients de l'erreur par rapport aux poids et aux biais, chacun décrit par les formules (26) et (27). Ils s'écrivent

$$\nabla_{\mathbf{w}^{(i)}} E = \nabla_{\mathbf{e}^{(i)}} E \nabla_{\mathbf{w}^{(i)}} \mathbf{e}^{(i)} = \nabla_{\mathbf{e}^{(i)}} E \mathbf{s}^{(i-1)\top}$$

et

$$\nabla_{\mathbf{b}^{(i)}} E = \nabla_{\mathbf{e}^{(i)}} E \nabla_{\mathbf{b}^{(i)}} \mathbf{e}^{(i)} = \nabla_{\mathbf{e}^{(i)}} E.$$

### 5.3 Calcul de la matrice hessienne dans une direction $\mathbf{v}$

Au lieu de parler de la dérivée seconde, il est maintenant question de trouver la hessienne de la fonction d'erreur. C'est celle-ci qui permet de donner une idée du comportement de la fonction d'erreur au point d'évaluation.

Par contre, la plupart des algorithmes d'optimisation s'intéressent plutôt à la hessienne dans une direction particulière  $v$  et non à toute la matrice hessienne. En d'autres termes, au lieu de calculer la hessienne  $\mathbf{H}$ , il est plus efficace de calculer directement le produit hessienne-vecteur  $\mathbf{H}\mathbf{v}$  puisque la matrice  $\mathbf{H}$  n'a pas à être gardée en mémoire.

À l'image de la version scalaire où le produit hessienne-vecteur de la fonction était composé des dérivées secondes de la fonction d'erreur par rapport aux poids et aux biais, la version vectorielle manipule directement les matrices de poids et de biais.

De la même façon que la version vectorielle de la propagation des données, la version vectorielle du calcul du produit hessienne-vecteur peut s'écrire

$$\begin{aligned}\mathcal{R}\{\mathbf{e}^{(i)}\} &= \mathcal{R}\{\mathbf{w}^{(i)}\mathbf{s}^{(i-1)} + \mathbf{b}^{(i)}\} \\ &= \mathbf{w}^{(i)}\mathcal{R}\{\mathbf{s}^{(i-1)}\} + \mathcal{R}\{\mathbf{w}^{(i)}\}\mathbf{s}^{(i-1)} + \mathcal{R}\{\mathbf{b}^{(i)}\} \\ &= \mathbf{w}^{(i)}\mathcal{R}\{\mathbf{s}^{(i-1)}\} + \mathbf{v}_w^{(i)}\mathbf{s}^{(i-1)} + \mathbf{v}_b^{(i)}\end{aligned}$$

et

$$\mathcal{R}\{\mathbf{s}^{(i)}\} = \mathbf{J}_{\mathbf{e}^{(i)}}(\mathbf{s}^{(i)}) \mathcal{R}\{\mathbf{e}^{(i)}\}.$$

Encore une fois, l'algorithme est initialisé avec  $\mathcal{R}\{\mathbf{s}^{(0)}\} = \mathbf{0}_{n \times 1}$ .

La méthode reste la même pour la rétropropagation. Par contre, il faut savoir ce que représente la dérivée seconde de la sortie d'un neurone par rapport à son entrée, c'est-à-dire de trouver  $\mathcal{R}\{\mathbf{J}_{\mathbf{e}^{(i)}}(\mathbf{s}^{(i)})\}$ . Or, comme la fonction d'activation n'a pas une image à une seule dimension, la matrice résultante serait en fait un tenseur (une matrice à plus de deux indices). Pour simplifier les choses, il est possible de se concentrer sur la deuxième forme de la formule (39) qui fait intervenir un vecteur multiplié à un autre vecteur de même dimension de façon terme-à-terme.

Dans cette notation, le  $j$ -ième élément du vecteur de la dérivée seconde est donné par

$$\text{diag}\left(\mathcal{R}\left\{\mathbf{J}_{\mathbf{e}^{(i)}}(\mathbf{s}^{(i)})\right\}\right)_j = s_j^{(i)}(1 - s_j^{(i)})(1 - 2s_j^{(i)})\mathcal{R}\left\{e_j^{(i)}\right\}.$$

Avec cette information, les formules (30), (31) et (32) deviennent respectivement

$$\begin{aligned}\mathcal{R}\{\nabla_{\mathbf{y}}E\} &= \mathcal{R}\{\mathbf{y} - \mathbf{c}\} \\ &= \mathcal{R}\{\mathbf{y}\}\end{aligned}\tag{40}$$

$$\begin{aligned}
\mathcal{R}\{\nabla_{\mathbf{e}^{(i)}}\} &= \mathcal{R}\left\{\mathbf{J}_{\mathbf{e}^{(i)}}\left(\mathbf{s}^{(i)}\right) \nabla_{\mathbf{s}^{(i)}} E\right\} \\
&= \mathcal{R}\left\{\mathbf{J}_{\mathbf{e}^{(i)}}\left(\mathbf{s}^{(i)}\right)\right\} \nabla_{\mathbf{s}^{(i)}} E + \mathbf{J}_{\mathbf{e}^{(i)}}\left(\mathbf{s}^{(i)}\right) \mathcal{R}\{\nabla_{\mathbf{s}^{(i)}} E\} \\
&= \text{diag}\left(\mathcal{R}\left\{\mathbf{J}_{\mathbf{e}^{(i)}}\left(\mathbf{s}^{(i)}\right)\right\}\right) \odot \nabla_{\mathbf{s}^{(i)}} E + \mathbf{J}_{\mathbf{e}^{(i)}}\left(\mathbf{s}^{(i)}\right) \mathcal{R}\{\nabla_{\mathbf{s}^{(i)}} E\}
\end{aligned}$$

$$\begin{aligned}
\mathcal{R}\{\nabla_{\mathbf{s}^{(i)}} E\} &= \mathcal{R}\left\{\left(\mathbf{w}^{(i+1)}\right)^\top \nabla_{\mathbf{e}^{(i+1)}} E\right\} \\
&= \left(\mathbf{w}^{(i+1)}\right)^\top \mathcal{R}\{\nabla_{\mathbf{e}^{(i+1)}} E\} + \mathcal{R}\left\{\left(\mathbf{w}^{(i+1)}\right)^\top\right\} \nabla_{\mathbf{e}^{(i+1)}} E \\
&= \left(\mathbf{w}^{(i+1)}\right)^\top \mathcal{R}\{\nabla_{\mathbf{e}^{(i+1)}} E\} + \left(\mathbf{v}_w^{(i+1)}\right)^\top \nabla_{\mathbf{e}^{(i+1)}} E
\end{aligned}$$

Encore une fois, à partir de ces résultats il est possible de reconstituer le vecteur du produit  $\mathbf{H}\mathbf{v}$  à partir de la concaténation des éléments des matrices, comme dans l'équation (35) :

$$\begin{aligned}
\mathcal{R}\{\nabla_{\mathbf{w}^{(i)}} E\} &= \mathcal{R}\left\{\nabla_{\mathbf{e}^{(i)}} E \left(\mathbf{s}^{(i-1)}\right)^\top\right\} \\
&= \nabla_{\mathbf{e}^{(i)}} E \mathcal{R}\left\{\left(\mathbf{s}^{(i-1)}\right)^\top\right\} + \mathcal{R}\{\nabla_{\mathbf{e}^{(i)}} E\} \left(\mathbf{s}^{(i-1)}\right)^\top
\end{aligned}$$

et des vecteurs

$$\mathcal{R}\{\nabla_{\mathbf{b}^{(i)}} E\} = \mathcal{R}\{\nabla_{\mathbf{e}^{(i)}} E\}$$

qui sont dérivés des formules (33) et (34).

## 6 Apprentissage d'un réseau de neurones

Comme il a été mentionné plus tôt, le but de la phase d'apprentissage est de trouver le vecteur de paramètres  $\boldsymbol{\theta}$  défini en (21) qui minimise la fonction d'erreur comme illustré à la formule (22).

Cette section porte sur quelques façons d'y arriver, chacune présentant leurs avantages et leurs inconvénients.

### 6.1 Descente de gradient

Une première méthode consiste à appliquer un algorithme de descente de gradient sur l'ensemble des couples  $(\mathbf{x}, \mathbf{c})$ . Cette façon de procéder consiste à trouver le gradient de l'erreur de chacun des couples et d'en faire la moyenne. Comme le gradient est un vecteur pointant vers la direction de plus forte croissance de la fonction d'erreur, déplacer les paramètres dans la direction opposée au gradient fait diminuer l'erreur.

Dans cette optique, les paramètres seront mis-à-jour selon l'équation  $\boldsymbol{\theta} := \boldsymbol{\theta} - \alpha (\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta})) = \alpha \left(\frac{1}{N} \sum_{t=1}^N \nabla E_t(\boldsymbol{\theta})\right)$  où  $\alpha$  est appelé le taux d'apprentissage

et contrôle la vitesse d'apprentissage. Un taux trop élevé pourrait nuire à l'algorithme qui sauterait par-dessus le déplacement idéal.

Selon cette méthode, l'algorithme d'apprentissage ressemble à l'algorithme 1. Dans cet algorithme, le critère d'arrêt peut être en lien avec un seuil d'erreur à atteindre ou un nombre limite d'itérations.

**Données** : Une base de données  $\mathbf{B}$  où  $\mathbf{B}_t$  est le  $t$ -ième couple  $(\mathbf{x}, \mathbf{c})$  et dont la taille est  $N$ .

**Entrées** : Une fonction `CritereArret()` qui renvoie vrai quand l'apprentissage est terminé.  
Un réseau de neurones muni d'une collection de paramètres  $\theta$ .  
Un hyper-paramètre  $\alpha$  donnant la vitesse d'apprentissage de l'algorithme.

**tant que** `CritereArret( $\theta$ )` **faire**

**pour chaque** *couple*  $(\mathbf{x}, \mathbf{c})$  **de**  $\mathbf{B}$  **faire**

$\mathbf{y} \leftarrow \text{Propagation}(\mathbf{x})$  ;                   /\* La prédiction du réseau \*/

$\mathbf{g}_t \leftarrow \text{Retropropagation}(\mathbf{y} - \mathbf{c})$  ;       /\* Le  $t$ -ième gradient \*/

$\mathbf{g} \leftarrow \frac{1}{N} \sum_{t=1}^N \mathbf{g}_t$  ;                               /\* Le gradient moyen \*/

$\theta \leftarrow \theta - \alpha \mathbf{g}$  ;                               /\* Mise à jour des paramètres \*/

$E \leftarrow \frac{1}{N} \sum_{t=1}^N \|y_t - c_t\|_2^2$  ;       /\* Erreur moyenne de prédiction \*/

**Sorties** : Un réseau de neurones ayant des paramètres  $\theta$  optimisés.

**Algorithme 1:** Algorithme de la descente de gradient pour l'apprentissage d'un réseau de neurones.

## 6.2 Descente de gradient par lots

Une autre méthode, l'apprentissage par lots<sup>3</sup> tente d'améliorer la vitesse d'apprentissage en utilisant un certain nombre de couples  $(\mathbf{x}, \mathbf{c})$  par lot plutôt que l'ensemble des éléments de la base de données.

Puisque les calculs ne se font pas sur l'ensemble des données, mais seulement sur un échantillon, l'estimation de la moyenne des gradients reste sans biais alors que la variance de l'estimation de la moyenne augmente. Or, il peut être plus efficace de procéder ainsi pour avancer plus rapidement dans l'espace solution. L'algorithme 2 présente le fonctionnement de cette méthode.

Une façon d'accélérer les choses, surtout dans l'optique de l'utilisation de bibliothèques d'algèbre linéaire optimisées, est de passer à la propagation tous les vecteurs d'entrée d'un lot à la fois, sous forme de matrice. De cette façon, les calculs peuvent tous se faire qu'en une seule propagation et rétropropagation.

Dans ce cas, un lot est un couple  $(\mathbf{X}, \mathbf{C})$  où  $\mathbf{X}$  et  $\mathbf{C}$  sont des matrices de dimensions  $n \times l$  composées de  $l$  vecteurs  $\mathbf{x}$  et  $\mathbf{c}$  respectivement.

À partir de ceci, chaque vecteur colonne des formules de la section 5 devient une matrice à  $l$  colonnes pour accommoder l'entrée et la cible de dimension  $l$ .

<sup>3</sup>De l'anglais *mini-batch training*.



est appelée l'apprentissage en ligne<sup>4</sup> lorsqu'elle se fait dans l'ordre des données ou apprentissage par descente de gradient stochastique lorsque le couple en entrée est choisi au hasard.

## 6.4 Méthodes d'optimisation de second ordre

Une autre possibilité est d'utiliser des algorithmes d'optimisation se servant d'informations sur la courbure de la fonction d'erreur, c'est-à-dire sur la matrice hessienne de  $E$ . Pour aider ces algorithmes, une méthode semblable à la rétro-propagation a été exposée dans cet article pour trouver le produit de la matrice hessienne et d'un vecteur. Ce produit est souvent utilisé dans ces techniques pour observer la courbure de la fonction d'erreur dans une direction donnée, ce qui ne nécessite pas le calcul de la hessienne en entier.

Bien entendu, il est possible de combiner ces méthodes avec la philosophie des lots afin d'accélérer l'apprentissage.

## Références

- [Bis06] C. M. BISHOP : *Pattern recognition and machine learning*. Springer New York, 2006.
- [Mø192] M. MØLLER : Supervised learning on large redundant training sets. *In Neural Networks for Signal Processing [1992] II., Proceedings of the 1992 IEEE-SP Workshop*, pages 79–89. IEEE, 1992.
- [Pea94] B. A. PEARLMUTTER : Fast exact multiplication by the hessian. *Neural Computation*, 6(1):147–160, 1994.
- [Roj96] R. ROJAS : *Neural networks : a systematic introduction*. Springer, 1996.
- [SB12] A. SALVAIL-BÉRARD : Les arbres de décision hybrides. *Cahier de Mathématique de l'Université de Sherbrooke*, 2:34–58, 2012.

ADAM SALVAIL-BÉRARD  
DÉPARTEMENT D'INFORMATIQUE, UNIVERSITÉ DE SHERBROOKE  
Courriel: Adam.Salvail-Berard@USherbrooke.ca

---

<sup>4</sup>De l'anglais *online training*.