

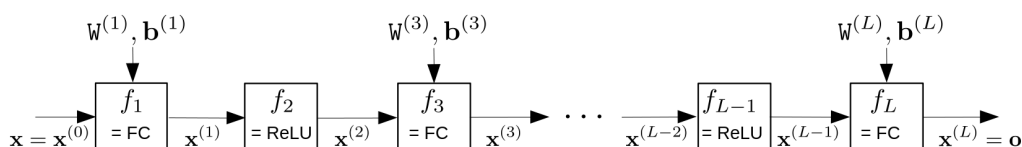
TD : Implémentation d'un réseau de neurones de type Perceptron multicouche pour un problème de classification

Guillaume Bourmaud

09/02/23

1 Introduction

Un réseau de neurones de type Perceptron multicouche (MLP) est composé d'une succession de couches, alternant une transformation affine (FC) et une fonction non-linéaire (historiquement la fonction sigmoïde mais ici nous utiliserons la fonction ReLU : $z = \max(0, x)$). Un tel réseau est illustré graphiquement ci-après :



Remarquons que selon les notations utilisées dans le schéma précédent, le nombre de couches L est nécessairement impair et supérieur ou égal à 3 (car les couches d'entrée et de sortie sont des transformations affines). Ainsi un MLP s'écrit comme une composition de fonctions :

$$\mathbf{o} = \text{MLP} \left(\mathbf{x}; \left\{ \mathbf{w}^{(2j-1)}, \mathbf{b}^{(2j-1)} \right\}_{j=1, \dots, \lceil L/2 \rceil} \right) \quad (1)$$

$$= \text{FC} \left(\text{ReLU} \left(\dots \text{FC} \left(\text{ReLU} \left(\text{FC} \left(\mathbf{x}; \mathbf{w}^{(1)}, \mathbf{b}^{(1)} \right) \right); \mathbf{w}^{(3)}, \mathbf{b}^{(3)} \right) \dots \right); \mathbf{w}^{(L)}, \mathbf{b}^{(L)} \right). \quad (2)$$

Nous considérons le cas où les vecteurs $(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(L-1)})$ sont tous de taille H et sont des vecteurs ligne, donc de taille $1 \times H$. Le MLP que nous considérons a donc deux hyper-paramètres : L (le nombre de couches) et H . Nous nous intéressons ici à un problème de classification à C classes, donc $\mathbf{x}^{(L)} = \mathbf{o}$ est de taille $1 \times C$. L'espace de départ est de dimension D donc $\mathbf{x} = \mathbf{x}^{(0)}$ est de taille $1 \times D$. La fonction de coût utilisée est l'entropie croisée Multinoulli :

$$l(y, \mathbf{o}) = -\ln(\mathbf{p}_y) \quad (3)$$

où $\mathbf{p} = \text{softmax}(\mathbf{o})$ est un vecteur de probabilités (scores positifs qui se somment à 1) de taille $1 \times C$ et $y \in \{0, \dots, C-1\}$ est l'étiquette (donc un scalaire 1×1). Par exemple, si $y = 2$, alors la notation \mathbf{p}_y correspond à récupérer le 3ème élément du vecteur \mathbf{p} . Minimiser la fonction de coût revient donc à maximiser le score de probabilité \mathbf{p}_y . La fonction softmax est définie de la manière suivante :

$$\mathbf{p}_i = \frac{\exp(\mathbf{o}_i)}{\sum_{j=0}^{C-1} \exp(\mathbf{o}_j)}. \quad (4)$$

Les paramètres du réseau seront optimisés en minimisant la fonction de coût suivante sur la base d'apprentissage $\{\mathbf{x}^{[i]}, y^{[i]}\}_{i=0 \dots N-1}$:

$$\arg \min_{\left\{ \mathbf{w}^{(2j-1)}, \mathbf{b}^{(2j-1)} \right\}_{j=1, \dots, \lceil L/2 \rceil}} \frac{1}{N} \sum_{i=0}^{N-1} l \left(y^{[i]}, \text{MLP} \left(\mathbf{x}^{[i]}; \left\{ \mathbf{w}^{(2j-1)}, \mathbf{b}^{(2j-1)} \right\}_{j=1, \dots, \lceil L/2 \rceil} \right) \right) \quad (5)$$

Nous utiliserons pour cela une méthode de descente de gradient. Comme vu en cours, le gradient, c'est-à-dire la dérivée de la fonction de coût par rapport aux paramètres du réseau, est obtenue en utilisant le théorème de dérivation des fonctions composées.

2 Théorème de dérivation d'une fonction composée

Prenons le cas d'une fonction composée de deux fonctions ($f : \mathbb{R}^D \rightarrow \mathbb{R}^H$ et $g : \mathbb{R}^H \rightarrow \mathbb{R}$):

$$a = g(f(\mathbf{x})). \quad (6)$$

En définissant $\mathbf{z} = f(\mathbf{x})$, le théorème de dérivation d'une fonction composée nous indique que la dérivée de $a \in \mathbb{R}$ par rapport au i -ème élément du vecteur \mathbf{x} s'exprime de la manière suivante :

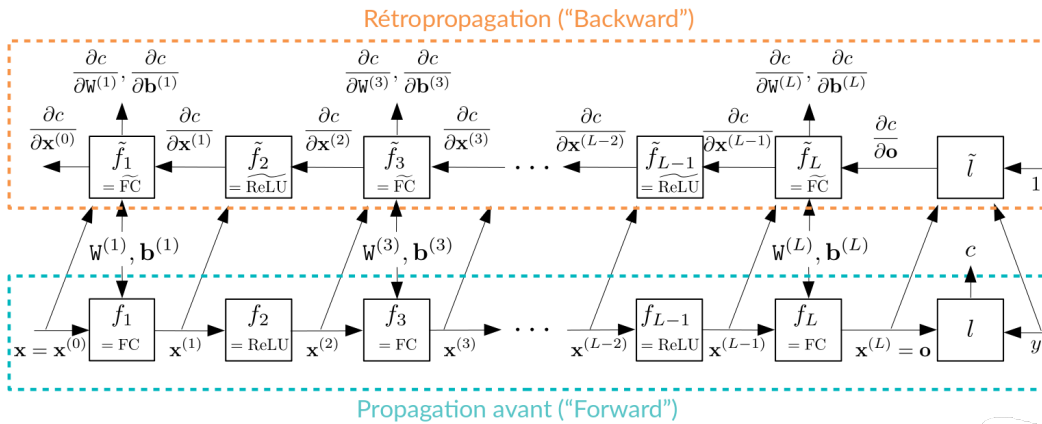
$$\frac{\partial a}{\partial \mathbf{x}_i} = \sum_{j=0}^{H-1} \frac{\partial a}{\partial \mathbf{z}_j} \bigg|_{\mathbf{y}=f(\mathbf{x})} \frac{\partial \mathbf{z}_j}{\partial \mathbf{x}_i}. \quad (7)$$

Ainsi nous pouvons introduire la notation suivante :

$$\frac{\partial a}{\partial \mathbf{x}} = \tilde{f} \left(\frac{\partial a}{\partial \mathbf{z}} \bigg|_{\mathbf{z}=f(\mathbf{x})}, \mathbf{x} \right), \quad (8)$$

où \tilde{f} est la fonction qui implémente l'équation (7), $\frac{\partial a}{\partial \mathbf{x}} = \left[\frac{\partial a}{\partial \mathbf{x}_0}, \frac{\partial a}{\partial \mathbf{x}_1}, \dots, \frac{\partial a}{\partial \mathbf{x}_{D-1}} \right]$ et $\frac{\partial a}{\partial \mathbf{z}} \bigg|_{\mathbf{z}=f(\mathbf{x})} = \left[\frac{\partial a}{\partial \mathbf{z}_0}, \frac{\partial a}{\partial \mathbf{z}_1}, \dots, \frac{\partial a}{\partial \mathbf{z}_{H-1}} \right]$.

En appliquant ce théorème à un MLP dans le but de calculer la dérivée de la fonction de coût par rapport aux paramètres du MLP, nous obtenons l'étape de rétropropagation du gradient, sous la forme d'un **graphe de calcul**, illustrée dans la figure ci-après.



Travail : écrire la taille de chacune des variables présentes dans le graphe de calcul précédent

Ainsi, pour mettre en œuvre une méthode de descente de gradient dans le but d'optimiser les paramètres d'un MLP, nous aurons uniquement besoin de connaître les expressions des fonctions $\tilde{l}(1, \mathbf{o})$, $\widetilde{\text{FC}} \left(\frac{\partial l}{\partial \mathbf{x}^{(i)}} \bigg|_{\mathbf{x}^{(i)}=\text{FC}(\mathbf{x}^{(i-1)})}, \mathbf{x}^{(i-1)}, \mathbf{W}^{(i)}, \mathbf{b}^{(i)} \right)$ et $\widetilde{\text{ReLU}} \left(\frac{\partial l}{\partial \mathbf{x}^{(i)}} \bigg|_{\mathbf{x}^{(i)}=\text{ReLU}(\mathbf{x}^{(i-1)})}, \mathbf{x}^{(i-1)} \right)$. C'est le travail qui sera réalisé dans la suite de ce TD. Les expressions obtenues seront codées ultérieurement lors d'un TP en Python utilisant la bibliothèque Numpy.

3 Notations

Dans le but de paralléliser les calculs, et ainsi d'obtenir une implémentation en Python efficace, nous considérons le cas où les N vecteurs de taille D $\{\mathbf{x}^{[i]}\}_{i=0 \dots N-1}$ de la base d'apprentissage sont organisés en une matrice :

$$\underbrace{\mathbf{X}}_{N \times D} = \begin{bmatrix} \mathbf{X}_{00} & \mathbf{X}_{01} & \dots & \mathbf{X}_{0(D-1)} \\ \mathbf{X}_{10} & \mathbf{X}_{11} & \dots & \mathbf{X}_{1(D-1)} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{X}_{(N-1)0} & \mathbf{X}_{(N-1)1} & \dots & \mathbf{X}_{(N-1)(D-1)} \end{bmatrix} \quad (9)$$

Les étiquettes $\{y^{[i]}\}_{i=0 \dots N-1}$ sont organisées sous la forme d'un vecteur :

$$\mathbf{y} = \begin{bmatrix} y_0 \\ y_1 \\ \vdots \\ y_{(N-1)} \end{bmatrix}^\top \quad (10)$$

Ainsi le MLP s'appliquera directement sur cette matrice \mathbf{X} et produira une matrice de scores \mathbf{O} :

$$\mathbf{O} = \text{MLP}(\mathbf{X}; \boldsymbol{\theta}) = \text{FC} \left(\text{ReLU} \left(\dots \text{FC} \left(\text{ReLU} \left(\text{FC} \left(\mathbf{X}; \mathbf{w}^{(1)}, \mathbf{b}^{(1)} \right) \right); \mathbf{w}^{(3)}, \mathbf{b}^{(3)} \right) \dots \right); \mathbf{w}^{(L)}, \mathbf{b}^{(L)} \right) \quad (11)$$

où

$$\underbrace{\mathbf{O}}_{N \times C} = \begin{bmatrix} \mathbf{O}_{00} & \mathbf{O}_{01} & \dots & \mathbf{O}_{0(C-1)} \\ \mathbf{O}_{10} & \mathbf{O}_{11} & \dots & \mathbf{O}_{1(C-1)} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{O}_{(N-1)0} & \mathbf{O}_{(N-1)1} & \dots & \mathbf{O}_{(N-1)(C-1)} \end{bmatrix} \quad (12)$$

Après softmax, les vecteurs de scores normalisés peuvent également être rangés dans une matrice :

$$\underbrace{\mathbf{P}}_{N \times C} = \begin{bmatrix} \mathbf{P}_{00} & \mathbf{P}_{01} & \dots & \mathbf{P}_{0(C-1)} \\ \mathbf{P}_{10} & \mathbf{P}_{11} & \dots & \mathbf{P}_{1(C-1)} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{P}_{(N-1)0} & \mathbf{P}_{(N-1)1} & \dots & \mathbf{P}_{(N-1)(C-1)} \end{bmatrix}. \quad (13)$$

Par définition, les matrices $(\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(L-1)})$ sont donc de taille $N \times H$.

4 Obtention des expressions de $\widetilde{\text{FC}}$, $\widetilde{\text{ReLU}}$ et \tilde{l}

4.1 Dérivée transformation affine : $\widetilde{\text{FC}}$

De manière générale, en utilisant les notations précédemment introduites, une transformation affine peut s'écrire :

$$z_{kl} = \sum_{j=0}^{D-1} x_{kj} w_{jl} + b_l \quad (14)$$

$$\text{où } \underbrace{\mathbf{Z}}_{N \times H} = \begin{bmatrix} \mathbf{Z}_{00} & \mathbf{Z}_{01} & \dots & \mathbf{Z}_{0(H-1)} \\ \mathbf{Z}_{10} & \mathbf{Z}_{11} & \dots & \mathbf{Z}_{1(H-1)} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{Z}_{(N-1)0} & \mathbf{Z}_{(N-1)1} & \dots & \mathbf{Z}_{(N-1)(H-1)} \end{bmatrix}, \underbrace{\mathbf{W}}_{D \times H} = \begin{bmatrix} \mathbf{W}_{00} & \mathbf{W}_{01} & \dots & \mathbf{W}_{0(H-1)} \\ \mathbf{W}_{10} & \mathbf{W}_{11} & \dots & \mathbf{W}_{1(H-1)} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{W}_{(D-1)0} & \mathbf{W}_{(D-1)1} & \dots & \mathbf{W}_{(D-1)(H-1)} \end{bmatrix}, \underbrace{\mathbf{b}}_{1 \times H} =$$

$$\begin{bmatrix} \mathbf{b}_0 \\ \mathbf{b}_1 \\ \vdots \\ \mathbf{b}_{(H-1)} \end{bmatrix}^\top$$

. En utilisant le produit matrice/vecteur, l'équation devient :

$$\mathbf{Z} = \text{FC}(\mathbf{X}, \mathbf{W}, \mathbf{b}) = \mathbf{X}\mathbf{W} + \mathbf{1}_N \mathbf{b} \quad (15)$$

où $\underbrace{\mathbf{1}_N}_{N \times 1} = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}$. La sortie \mathbf{Z} de cette transformation affine est passée dans une fonction g dont la sortie s est scalaire :

$$s = g(\mathbf{Z}) \quad (16)$$

En supposant le gradient $\underbrace{\partial s / \partial \mathbf{z}}_{N \times H} = \begin{bmatrix} \partial s / \partial \mathbf{z}_{00} & \partial s / \partial \mathbf{z}_{01} & \dots & \partial s / \partial \mathbf{z}_{0(H-1)} \\ \partial s / \partial \mathbf{z}_{10} & \partial s / \partial \mathbf{z}_{11} & \dots & \partial s / \partial \mathbf{z}_{1(H-1)} \\ \vdots & \vdots & \ddots & \vdots \\ \partial s / \partial \mathbf{z}_{(N-1)0} & \partial s / \partial \mathbf{z}_{(N-1)1} & \dots & \partial s / \partial \mathbf{z}_{(N-1)(H-1)} \end{bmatrix}$ connu, l'objectif est de

calculer la dérivée de s par rapport aux entrées \mathbf{X} et aux paramètres \mathbf{W} et \mathbf{b} .

Travail : dessiner le graphe de calcul correspondant aux deux équations précédentes, en faisant apparaître $s, g, \mathbf{Z}, \mathbf{FC}, \mathbf{X}, \mathbf{W}, \mathbf{b}, \tilde{g}, \partial s / \partial \mathbf{z}, \widetilde{\mathbf{FC}}, \partial s / \partial \mathbf{x}, \partial s / \partial \mathbf{w}$ et $\partial s / \partial \mathbf{b}$.

L'objectif de cette partie est d'obtenir des expressions permettant d'implémenter la fonction $\widetilde{\mathbf{FC}}$. Il faut donc calculer les expressions de $\partial s / \partial \mathbf{x}$, $\partial s / \partial \mathbf{w}$ et $\partial s / \partial \mathbf{b}$ en les simplifiant au maximum pour faire apparaître des opérations permettant une implémentation efficace en Python (par exemple des produits de matrices).

4.1.1 Calcul de $\partial s / \partial \mathbf{x}$

Le théorème de dérivation d'une fonction composée nous dit que :

$$\partial s / \partial \mathbf{x}_{ij} = \sum_{k=0}^{N-1} \sum_{l=0}^{H-1} \partial s / \partial \mathbf{z}_{kl} \cdot \partial \mathbf{z}_{kl} / \partial \mathbf{x}_{ij} \quad (17)$$

Travail :

- montrer que l'équation précédente se simplifie en $\partial s / \partial \mathbf{x}_{ij} = \sum_{l=0}^{H-1} \partial s / \partial \mathbf{z}_{il} \cdot \mathbf{W}_{jl}$

• sachant que $\underbrace{\partial s / \partial \mathbf{x}}_{N \times D} = \begin{bmatrix} \partial s / \partial \mathbf{x}_{00} & \partial s / \partial \mathbf{x}_{01} & \dots & \partial s / \partial \mathbf{x}_{0(D-1)} \\ \partial s / \partial \mathbf{x}_{10} & \partial s / \partial \mathbf{x}_{11} & \dots & \partial s / \partial \mathbf{x}_{1(D-1)} \\ \vdots & \vdots & \ddots & \vdots \\ \partial s / \partial \mathbf{x}_{(N-1)0} & \partial s / \partial \mathbf{x}_{(N-1)1} & \dots & \partial s / \partial \mathbf{x}_{(N-1)(D-1)} \end{bmatrix}$, montrer que $\underbrace{\partial s / \partial \mathbf{x}}_{N \times D} = \underbrace{\partial s / \partial \mathbf{z}}_{N \times H} \underbrace{\mathbf{W}^\top}_{H \times D}$

Ainsi $\partial s / \partial \mathbf{x}$ pourra être calculé lors de l'implémentation en Python par un simple appel à une fonction réalisant un produit matriciel.

4.1.2 Calcul de $\partial s / \partial \mathbf{w}$

Le théorème de dérivation d'une fonction composée nous dit que :

$$\partial s / \partial \mathbf{w}_{ij} = \sum_{k=0}^{N-1} \sum_{l=0}^{H-1} \partial s / \partial \mathbf{z}_{kl} \cdot \partial \mathbf{z}_{kl} / \partial \mathbf{w}_{ij} \quad (18)$$

Travail :

- montrer que l'équation précédente se simplifie en $\partial s / \partial \mathbf{w}_{ij} = \sum_{k=0}^{N-1} \partial s / \partial \mathbf{z}_{kj} \cdot \mathbf{x}_{ki}$

• sachant que $\underbrace{\partial s / \partial \mathbf{w}}_{D \times H} = \begin{bmatrix} \partial s / \partial \mathbf{w}_{00} & \partial s / \partial \mathbf{w}_{01} & \dots & \partial s / \partial \mathbf{w}_{0(H-1)} \\ \partial s / \partial \mathbf{w}_{10} & \partial s / \partial \mathbf{w}_{11} & \dots & \partial s / \partial \mathbf{w}_{1(H-1)} \\ \vdots & \vdots & \ddots & \vdots \\ \partial s / \partial \mathbf{w}_{(D-1)0} & \partial s / \partial \mathbf{w}_{(D-1)1} & \dots & \partial s / \partial \mathbf{w}_{(D-1)(H-1)} \end{bmatrix}$ montrer que $\underbrace{\partial s / \partial \mathbf{w}}_{D \times H} = \underbrace{\mathbf{x}^\top}_{D \times N} \underbrace{\partial s / \partial \mathbf{z}}_{N \times H}$

Ainsi, comme $\partial s / \partial \mathbf{x}$, $\partial s / \partial \mathbf{w}$ pourra être calculé lors de l'implémentation en Python par un simple appel à une fonction réalisant un produit matriciel.

4.1.3 Calcul de $\partial s / \partial \mathbf{b}$

Travail :

- montrer que $\underbrace{\partial s / \partial \mathbf{b}}_{1 \times H} = \underbrace{\mathbf{1}_N^\top}_{1 \times N} \underbrace{\partial s / \partial \mathbf{z}}_{N \times H}$

Ainsi $\partial s / \partial \mathbf{b}$ pourra être calculé lors de l'implémentation en Python par un simple appel à une fonction réalisant un produit vecteur/matrice.

4.2 Dérivée ReLU : $\widetilde{\text{ReLU}}$

De manière générale, en utilisant les notations précédemment introduites, la fonction ReLU peut s'écrire :

$$Z_{kl} = \max(0, X_{kl}) \quad (19)$$

La sortie Z de cette transformation affine est passée dans une fonction g dont la sortie s est scalaire :

$$s = g(Z) \quad (20)$$

En supposant le gradient $\underbrace{\partial s / \partial Z}_{N \times D} = \begin{bmatrix} \partial s / \partial Z_{00} & \partial s / \partial Z_{01} & \dots & \partial s / \partial Z_{0(D-1)} \\ \partial s / \partial Z_{10} & \partial s / \partial Z_{11} & \dots & \partial s / \partial Z_{1(D-1)} \\ \vdots & \vdots & \ddots & \vdots \\ \partial s / \partial Z_{(N-1)0} & \partial s / \partial Z_{(N-1)1} & \dots & \partial s / \partial Z_{(N-1)(D-1)} \end{bmatrix}$ connu, l'objectif est de calculer la dérivée de s par rapport aux entrées X .

Travail : dessiner le graphe de calcul correspondant aux deux équations précédentes, en faisant apparaître $s, g, Z, \text{ReLU}, X, \tilde{g}, \partial s / \partial Z, \widetilde{\text{ReLU}}$ et $\partial s / \partial X$.

- **montrer que** $\partial s / \partial X_{kl} = \begin{cases} \partial s / \partial Z_{kl} & \text{si } X_{kl} > 0 \\ 0 & \text{si } X_{kl} < 0 \end{cases}$
- **la dérivée est-elle définie en 0 ?**

4.3 Dérivée de la fonction de coût : \tilde{l}

De manière générale, en utilisant les notations précédemment introduites, la fonction de coût l peut s'écrire :

$$r = \frac{1}{N} \sum_{i=0}^{N-1} -\ln \left(\frac{\exp(0_{iy_i})}{\sum_{j=0}^{C-1} \exp(0_{ij})} \right) \quad (21)$$

Travail : dessiner le graphe de calcul correspondant à l'équation précédente, en faisant apparaître $r, l, 0, \tilde{l}, \partial r / \partial 0$.

Travail :

- **montrer que** $\partial r / \partial 0_{ik} = \begin{cases} \frac{1}{N} (P_{ik} - 1) & \text{si } k = y_i \\ \frac{P_{ik}}{N} & \text{sinon} \end{cases}$

Les différentes expressions obtenues au cours de ce TD seront utilisées au prochain TP pour implémenter en Python (à l'aide de la bibliothèque Numpy) un MLP.