

Chapitre 6.2 - Le réseau de neurones

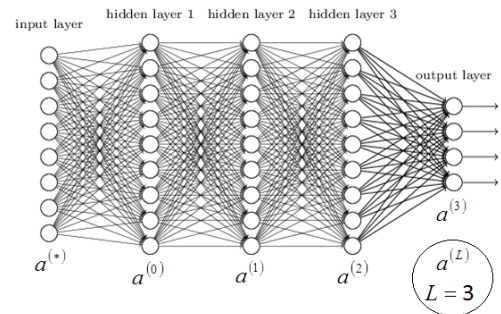
LES CONCEPTS DE BASE	2
LE RÉSEAU DE NEURONES	2
LE SIGNAL D'UN NEURONE.....	2
LES POIDS D'UN NEURONE.....	3
LE BIAIS D'UN NEURONE.....	3
LES FONCTIONS	4
LA FONCTION D'AGRÉGATION.....	4
<i>La différentielle de la fonction d'agrégation.....</i>	4
<i>La fonction de la somme pondérée.....</i>	5
<i>La fonction de la somme quadratique.....</i>	6
LA FONCTION D'ACTIVATION	7
<i>La différentielle de la fonction d'activation</i>	7
<i>La fonction sigmoïde.....</i>	8
<i>La fonction tanh.....</i>	9
<i>La fonction ReLU.....</i>	10
<i>La fonction Leaky ReLU.....</i>	11
<i>La fonction Randomized Leaky ReLU (2015).....</i>	11
<i>La fonction ELU.....</i>	11
<i>La fonction swish.....</i>	12
<i>La fonction softmax.....</i>	15
<i>La fonction softmax numériquement stable.....</i>	17
LA FONCTION D'ERREUR.....	18
<i>La différentielle de la fonction d'erreur.....</i>	18
<i>La fonction d'erreur linéaire.....</i>	19
<i>La fonction d'erreur quadratique.....</i>	19
<i>La fonction d'erreur cross entropy.....</i>	19
L'APPRENTISSAGE.....	20
LE GRADIENT DE LA FONCTION D'ERREUR.....	20
L'ALGORITHME DE L'APPRENTISSAGE	21
LA PROPAGATION DU GRADIENT	22
LES ÉQUATIONS DE LA PROPAGATION DU GRADIENT.....	25
L'ÉQUATION DE LA PROPAGATION DE L'ERREUR DE LA FONCTION D'ERREUR.....	25
L'ÉQUATION DE LA PROPAGATION DE L'ERREUR DE LA FONCTION D'ACTIVATION	25
L'ÉQUATION DE LA PROPAGATION DE L'ERREUR DE LA FONCTION D'ACTIVATION À PLUSIEURS NEURONES	26
L'ÉQUATION DE LA PROPAGATION DE L'ERREUR DE LA FONCTION D'AGRÉGATION.....	27
LE TABLEAU RÉCAPITULATIF DES ÉQUATIONS DE LA PROPAGATION DU GRADIENT	29
APPENDICE	30
LES NOTATIONS.....	30

Les concepts de base

Le réseau de neurones

Un réseau de neurones est une structure mathématique ordonnée en séquence de couches de donnée permettant de générer à partir d'un vecteur d'entrée $a^{(*)}$ un vecteur de sortie $a^{(L)}$. Le but du réseau sera de faire cheminer un vecteur $a^{(*)}$ par l'action de l'activation au travers l'ensemble des couches du réseau afin d'extraire une caractéristique au vecteur d'entrée initial sous la forme d'un vecteur de sortie $a^{(L)}$. Le réseau de neurone joue ainsi le rôle de fonction mathématique.

L'activation d'un réseau se fait en séquence de couche. Dans la notation retenue, nous utiliserons la séquence suivante pour un réseau contenant M couches de neurones :



Exemple de réseau pleinement connecté (full connected network)

Couche #1		Couche #2		...	Couche #m		...	Dernière couche #M $M = L + 1$	
$k = 0$		$k = 1$			$k = m - 1$			$k = L$	
Entrée	Sortie	Entrée	Sortie		Entrée	Sortie		Entrée	Sortie
$a^{(*)}$	$a^{(0)}$	$a^{(0)}$	$a^{(1)}$		$a^{(k-1)}$	$a^{(k)}$		$a^{(L-1)}$	$a^{(L)}$

Le signal d'un neurone

Dans un réseau de neurone, on utilise l'expression « signal » pour représenter les composantes d'un vecteur. On distingue deux types de signaux pour une couche k : le signal d'entrée $a^{(k-1)}$ et le signal de sortie $a^{(k)}$ (l'activation d'une couche). Idéalement, les composantes $a_i^{(k)}$ doivent être situées entre -1 et 1 pour augmenter la précision du réseau :

- $a_i^{(*)}$: Entrée (input) du neurone i du réseau de neurone où $i \in [0, N^{(*)} - 1]$. Il appartient à une donnée initiale (vecteur d'entrée) ou à la sortie i d'un autre réseau de neurones.
- $a_i^{(0)}$: Sortie (output) du neurone i de la 1^{re} couche de neurones du réseau ($k = 0$) avec $i \in [0, N^{(0)} - 1]$. Il sera l'entrée pour la 2^e couche de neurones du réseau ($k = 1$)
- $a_i^{(k)}$: Sortie (output) du neurone i de la k ^e couche de neurones du réseau où $i \in [0, N^{(k)} - 1]$. Il sera l'entrée pour la $k+1$ ^e couche du réseau.
- $a_i^{(L)}$: Sortie (output) du neurone i du réseau de neurones où $i \in [0, N^{(L)} - 1]$. Il appartient à la dernière couche du réseau ($k = L$).



<https://www.cgtrader.com/3d-models/character/anatomy/nerve-cell-anatomy-in-details>

Le signal d'un neurone d'un réseau artificiel a été modélisé sur le concept biologique du neurone que l'on retrouve dans le cerveau.

Les poids d'un neurone

Dans un réseau de neurones, les poids $w^{(k)}$ sont utilisés pour représenter une connectivité entre les signaux d'entrée $a^{(k-1)}$ à une couche k afin de générer un signal de sortie $a^{(k)}$ à la couche k . Cette connectivité permet d'extraire des caractéristiques entre les données $a^{(k-1)}$ et les valeurs des poids $w^{(k)}$ aura pour but de les identifier.

Le poids $w_{ij}^{(k)}$ établira la connexion entre le neurone d'entrée $a_j^{(k-1)}$ et le neurone de sortie $a_i^{(k)}$. Il déterminera l'importance du neurone $a_j^{(k-1)}$ dans l'émission du signal $a_i^{(k)}$. La valeur de $w_{ij}^{(k)}$ est un nombre réel pouvant être positif ou négatif.

Dans une représentation en algèbre linéaire

$$y = mx + b,$$

les poids $w_{ij}^{(k)}$ correspondent à la variable m dans l'équation, le neurone d'entrée $a_j^{(k-1)}$ représente la variable x et le neurone de sortie $a_i^{(k)}$ représente la variable y . Le poids $w^{(*)}$ n'existe pas.

On distinguera deux notations pour les poids d'un réseau :

Notation matricielle (notation ligne-colonne)	Notation tensorielle (notation in-out)
$w_{uv}^{(k)}$ (Composant de matrice en ligne u et colonne v)	$\tilde{w}_{vu}^{(k)}$ (Composant du tenseur en entrée v et sortie u)

Ainsi, en notation matricielle, la ligne de la matrice désigne l'ensemble des poids qui seront exploités pour évaluer le signal de sortie d'un neurone.

Le biais d'un neurone

Le biais $b^{(k)}$ (*bias*) est l'ensemble des poids constants applicable au calcul des signaux $a^{(k)}$ appartenant à une couche k .

Le biais $b_i^{(k)}$ est une constants lors du calcul d'un signal $a_i^{(k)}$ appartenant à une couche k . La valeur de $b_i^{(k)}$ est un nombre réel pouvant être positif ou négatif.

Dans une représentation en algèbre linéaire

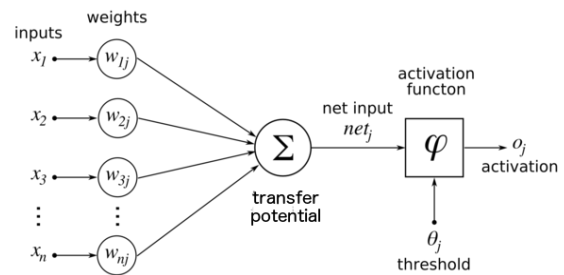
$$y = mx + b,$$

les biais $b^{(k)}$ correspondent à la variable b dans l'équation. Le biais $b^{(*)}$ n'existe pas.

Les fonctions

La fonction d'agrégation

L'agrégation est l'action de regrouper l'ensemble des neurones $a^{(k-1)}$ d'une couche $k-1$ afin de calculer un terme d'agrégation $z_i^{(k)}$ permettant de déterminer la valeur d'un neurone $a_i^{(k)}$ de sortie de la couche k . Il existe plusieurs fonctions d'agrégation. La fonction la plus utilisée est l'agrégation de la somme pondérée. Il est à noter que $z_i^{(*)}$ n'existe pas.



<https://medium.com/autonomous-agents/mathematical-foundation-for-activation-functions-in-artificial-neural-networks-a51c9dd7c089>

Schéma de n entrées pour déterminer l'agrégation d'un neurone sans l'usage de biais.

Mathématiquement, la fonction d'agrégation prendra la forme de l'équation suivante :

$$z_u^{(k)} = z_u^{(k)}(a_v^{(k-1)}, \tilde{w}_{vu}^{(k)}, b_u^{(k)}) \quad \text{avec } v \in [0, N^{(k-1)} - 1] \text{ et } u \in [0, N^{(k)} - 1]$$

où u : Indice de la sortie (*output*) de la fonction.

v : Indice d'entrée (*input*) de la fonction.

$z_u^{(k)}$: Agrégation du neurone u de la couche k .

$a_v^{(k-1)}$: Entrée du neurone v de la couche $k-1$.

$\tilde{w}_{vu}^{(k)}$: Poids du neurone d'entrée v sur le neurone de sortie u appartenant à la couche k .

$b_u^{(k)}$: Biais du neurone de sortie u appartenant à la couche k .

La différentielle de la fonction d'agrégation

L'expression de la différentielle de la fonction d'agrégation

$$z_u^{(k)} = z_u^{(k)}(a_v^{(k-1)}, \tilde{w}_{vu}^{(k)}, b_u^{(k)})$$

où

$$v \in [0 \dots N^{(k-1)} - 1] \quad \text{et} \quad u \in [0 \dots N^{(k)} - 1]$$

est égale à l'expression suivante :

$$dz_u^{(k)} = \sum_{v=0}^{N^{(k-1)}-1} \left(\frac{\partial z_u^{(k)}}{\partial a_v^{(k-1)}} da_v^{(k-1)} + \frac{\partial z_u^{(k)}}{\partial \tilde{w}_{vu}^{(k)}} d\tilde{w}_{vu}^{(k)} \right) + \frac{\partial z_u^{(k)}}{\partial b_u^{(k)}} db_u^{(k)}$$

(Différentielle de la fonction d'agrégation)

Preuve :

En construction ...

La fonction de la somme pondérée

La somme pondérée (*weighted sum*) est une fonction d'agrégation correspondant à une équation linéaire de type

$$y = mx + b$$

où x est la variable d'entrée, y est la variable de sortie, m est la constante linéaire et b est une constante.

La somme pondérée (*weighted sum*) $z_i^{(k)}$ d'un neurone i appartenant à la couche k correspond à l'expression suivante qui est déterminée à partir des entrées $a_j^{(k-1)}$ de la couche précédente ou $j \in [0, N^{(k-1)}]$:

$$z_i^{(k)} = \sum_{j=0}^{N^{(k-1)}-1} w_{ij}^{(k)} a_j^{(k-1)} + b_i^{(k)}$$

(Notation des poids en matrice ligne-colonne)

$$z_i^{(k)} = \sum_{j=0}^{N^{(k-1)}-1} \tilde{w}_{ji}^{(k)} a_j^{(k-1)} + b_i^{(k)}$$

(Notation des poids en tenseur *in-out*)

En notation matricielle, la somme pondérée correspond à l'équation suivante :

$$\begin{pmatrix} z_0^{(k)} \\ z_1^{(k)} \\ \dots \\ z_u^{(k)} \\ \dots \\ z_G^{(k)} \end{pmatrix} = \begin{pmatrix} w_{00}^{(k)} & w_{01}^{(k)} & \dots & w_{0v}^{(k)} & \dots & \dots & w_{0F}^{(k)} \\ w_{10}^{(k)} & w_{11}^{(k)} & \dots & w_{1v}^{(k)} & \dots & \dots & w_{1F}^{(k)} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ w_{u0}^{(k)} & w_{u1}^{(k)} & \dots & w_{uv}^{(k)} & \dots & \dots & w_{uF}^{(k)} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ w_{G0}^{(k)} & w_{G1}^{(k)} & \dots & w_{Gv}^{(k)} & \dots & \dots & w_{GF}^{(k)} \end{pmatrix} \begin{pmatrix} a_0^{(k-1)} \\ a_1^{(k-1)} \\ \dots \\ a_v^{(k-1)} \\ \dots \\ a_F^{(k-1)} \end{pmatrix} + \begin{pmatrix} b_0^{(k)} \\ b_1^{(k)} \\ \dots \\ b_u^{(k)} \\ \dots \\ b_G^{(k)} \end{pmatrix}$$

où $F = N^{(k-1)} - 1$: Indice du dernier neurone en entrée (indice de colonne de la matrice).

$G = N^{(k)} - 1$: Indice du dernier neurone en sortie (indice de ligne de la matrice).

Les dérivées partielles de la somme pondérée donnent les résultats suivants :

$$\frac{\partial z_i^{(k)}}{\partial w_{ij}^{(k)}} = a_j^{(k-1)}$$

(Dérivée par rapport au poids)

$$\frac{\partial z_i^{(k)}}{\partial a_j^{(k-1)}} = w_{ij}^{(k)}$$

(Dérivée par rapport à l'entrée)

$$\frac{\partial z_i^{(k)}}{\partial b_i^{(k)}} = 1$$

(Dérivée par rapport au biais)

La fonction de la somme quadratique

La somme quadratique (*quadratic sum*) est une fonction d'agrégation correspondant à une équation du 2^e degré de type

$$y = ax^2 + bx + c$$

où x est la variable d'entrée, y est la variable de sortie, a est la constante quadratique, b est la constante linéaire et c est une constante.

La somme quadratique¹ (*quadratic sum*) $z_i^{(k)}$ d'un neurone i appartenant à la couche k correspond à l'expression suivante qui est déterminée à partir des entrées $a_j^{(k-1)}$ ou $j \in [0, N^{(k-1)}]$ de la couche précédente :

$$z_i^{(k)} = \sum_{j=0}^{N^{(k-1)}-1} (w_{ij}^{1(k)} a_j^{(k-1)} + w_{ij}^{2(k)} a_j^{(k-1)} a_j^{(k-1)}) + b_i^{(k)} \quad z_i^{(k)} = \sum_{i=0}^{N^{(k-1)}-1} (\tilde{w}_{ji}^{1(k)} a_j^{(k-1)} + \tilde{w}_{ji}^{2(k)} a_j^{(k-1)} a_j^{(k-1)}) + b_i^{(k)}$$

(Notation des poids en matrice ligne-colonne)

(Notation des poids en tenseur *in-out*)

En notation matricielle, la somme pondérée correspond à l'équation suivante :

$$\begin{pmatrix} z_0^{(k)} \\ z_1^{(k)} \\ \dots \\ z_u^{(k)} \\ \dots \\ z_G^{(k)} \end{pmatrix} = \begin{pmatrix} w_{00}^{1(k)} & w_{01}^{1(k)} & \dots & w_{0v}^{1(k)} & \dots & \dots & w_{0F}^{1(k)} \\ w_{10}^{1(k)} & w_{11}^{1(k)} & \dots & w_{1v}^{1(k)} & \dots & \dots & w_{1F}^{1(k)} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ w_{u0}^{1(k)} & w_{u1}^{1(k)} & \dots & w_{uv}^{1(k)} & \dots & \dots & w_{uF}^{1(k)} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ w_{G0}^{1(k)} & w_{G1}^{1(k)} & \dots & w_{Gv}^{1(k)} & \dots & \dots & w_{GF}^{1(k)} \end{pmatrix} \begin{pmatrix} a_0^{(k-1)} \\ a_1^{(k-1)} \\ \dots \\ a_v^{(k-1)} \\ \dots \\ a_F^{(k-1)} \end{pmatrix} + \begin{pmatrix} w_{00}^{2(k)} & w_{01}^{2(k)} & \dots & w_{0v}^{2(k)} & \dots & \dots & w_{0F}^{2(k)} \\ w_{10}^{2(k)} & w_{11}^{2(k)} & \dots & w_{1v}^{2(k)} & \dots & \dots & w_{1F}^{2(k)} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ w_{u0}^{2(k)} & w_{u1}^{2(k)} & \dots & w_{uv}^{2(k)} & \dots & \dots & w_{uF}^{2(k)} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ w_{G0}^{2(k)} & w_{G1}^{2(k)} & \dots & w_{Gv}^{2(k)} & \dots & \dots & w_{GF}^{2(k)} \end{pmatrix} \begin{pmatrix} a_0^{(k-1)} * a_0^{(k-1)} \\ a_1^{(k-1)} * a_1^{(k-1)} \\ \dots \\ a_v^{(k-1)} * a_v^{(k-1)} \\ \dots \\ a_F^{(k-1)} * a_F^{(k-1)} \end{pmatrix} + \begin{pmatrix} b_0^{(k)} \\ b_1^{(k)} \\ \dots \\ b_u^{(k)} \\ \dots \\ b_G^{(k)} \end{pmatrix}$$

où $F = N^{(k-1)} - 1$: Indice du dernier neurone en entrée (indice de colonne de la matrice).

$G = N^{(k)} - 1$: Indice du dernier neurone en sortie (indice de ligne de la matrice).

Les dérivées partielles de la somme quadratique donnent les résultats suivants :

$$\frac{\partial z_i^{(k)}}{\partial w_{ij}^{1(k)}} = a_j^{(k-1)} \quad \frac{\partial z_i^{(k)}}{\partial w_{ij}^{2(k)}} = a_j^{(k-1)} a_j^{(k-1)} \quad \frac{\partial z_i^{(k)}}{\partial a_j^{(k-1)}} = w_{ij}^{1(k)} + 2a_j^{(k-1)} w_{ij}^{2(k)} \quad \frac{\partial z_i^{(k)}}{\partial b_i^{(k)}} = 1$$

(Dérivée par rapport au poids d'ordre 1)

(Dérivée par rapport au poids d'ordre 2)

(Dérivée par rapport à l'entrée)

(Dérivée par rapport au biais)

¹ L'efficacité de cette fonction est présentement en étude. Il en existe plusieurs autres formes équivalentes.

La fonction d'activation

L'activation est l'action de transformer le résultat de l'agrégation $z_i^{(k)}$ d'un neurone i appartenant à la couche k en une valeur située idéalement entre -1 et 1. Il existe plusieurs fonctions d'activation avec plusieurs types de valeurs de sortie. Ceci permet de donner un comportement non linéaire au réseau de neurones.

La fonction d'activation $f^{(k)}$ qui génère l'ensemble des neurones $a^{(k)}$ appartenant à la couche k s'applique sur les résultats de l'agrégation $z^{(k)}$ de la façon suivante :

$$a_i^{(k)} = f_i^{(k)}(z_i^{(k)}) \quad \text{et} \quad i \in [0, N^{(k)} - 1]$$

où $z_i^{(k)}$ est le paramètre de la fonction $f^{(k)}$. La fonction d'activation $f^{(*)}$ n'existe pas.

Il existe des fonctions d'activation dont l'expression dépend de la valeur de l'agrégation de tous les neurones (*full connected*). Elle peut alors prendre la forme suivante :

$$a_i^{full(k)} = f_i^{full(k)}(\bar{z}^{(k)}) = f_i^{full(k)}(z_0^{(k)}, z_1^{(k)}, \dots, z_i^{(k)}, \dots, z_{N^{(k)}-1}^{(k)})$$

La différentielle de la fonction d'activation

La différentielle de la fonction d'activation dépend du nombre de neurones en entrée de celle-ci. Puisque nous avons la fonction d'activation à un neurone

$$a_u^{(k)} = a_u^{(k)}(z_u^{(k)})$$

et la fonction d'activation à plusieurs neurones

$$a_u^{full(k)} = a_u^{full(k)}(\bar{z}^{(k)}) = a_u^{full(k)}(z_0^{(k)}, z_1^{(k)}, \dots, z_p^{(k)}, \dots, z_{N^{(k)}-1}^{(k)})$$

où

$$u \in [0 \dots N^{(k)} - 1],$$

nous appliquons la règle de la dérivée en chaîne et nous obtenons les expression suivante pour la différentielle de nos fonctions d'activations :

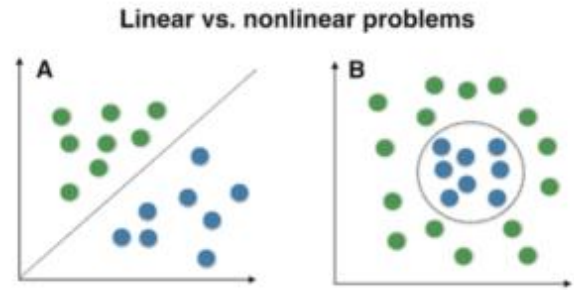
Différentielle de la fonction d'activation à un neurone	Différentielle de la fonction d'activation à plusieurs neurones ² (<i>full connected</i>)	Différentielle de la fonction d'activation avec paramètre β ³
$da_u^{(k)} = \frac{\partial a_u^{(k)}}{\partial z_u^{(k)}} dz_u^{(k)}$	$da_u^{full(k)} = \sum_{p=0}^{N^{(k)}-1} \frac{\partial a_u^{(k)}}{\partial z_p^{(k)}} dz_p^{(k)}$	$da_u^{(k)} = \frac{\partial a_u^{(k)}}{\partial z_u^{(k)}} dz_u^{(k)} + \frac{\partial a_u^{(k)}}{\partial \beta^{(k)}} d\beta^{(k)}$

Preuve :

En construction ...

² Voir la fonction *softmax*.

³ Voir la fonction *swish*.



<https://www.datarobot.com/wiki/neural-network/>

Illustration d'une classification à l'aide d'une fonction d'activation linéaire et non linéaire.

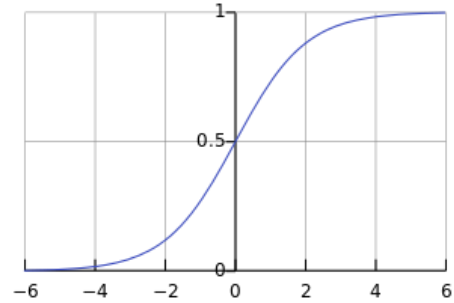
La fonction sigmoïde

La fonction sigmoïde est une fonction d'activation dont le résultat est situé dans l'intervalle $[0,1]$:

$$\text{Fonction : } \sigma(z) = \frac{1}{1 + e^{-z}}$$

$$\text{Dérivée : } \frac{d\sigma(z)}{dz} = \sigma(z)(1 - \sigma(z))$$

De nos jours, la fonction sigmoïde est évitée dans l'activation d'une couche de neurones en raison de la faible valeur de sa dérivée lorsque $|z| > 4$. Puisque le taux d'apprentissage dépend de la dérivée de la fonction d'activation, si $|z|$ est trop élevée alors la dérivée chutera vers zéro stoppant l'apprentissage.



<https://www.quora.com/Why-is-the-sigmoid-function-rarely-used-in-hidden-layers-recently>

Illustration des valeurs de la fonction sigmoïde.

Preuve :

Effectuons la dérivée de la fonction sigmoïde $\sigma(x)$ par rapport à la variable x :

$$\begin{aligned} \sigma' &= \frac{d\sigma}{dx} & \Rightarrow & \sigma' = \frac{d}{dx} \left(\frac{1}{1 + e^{-x}} \right) & (\sigma(x) &= \frac{1}{1 + e^{-x}}) \\ & & \Rightarrow & \sigma' = -\frac{1}{(1 + e^{-x})^2} \frac{d}{dx} (1 + e^{-x}) & \left(\frac{d}{dx} \left(\frac{1}{f} \right) = -\frac{1}{f^2} \frac{df}{dx} \right) \\ & & \Rightarrow & \sigma' = \frac{-1}{(1 + e^{-x})^2} e^{-x} \frac{d}{dx} (-x) & \left(\frac{d}{dx} (e^{-x}) = e^{-x} \frac{d}{dx} (-x) \right) \\ & & \Rightarrow & \sigma' = \frac{-e^{-x}}{(1 + e^{-x})^2} \frac{d}{dx} (-x) \\ & & \Rightarrow & \sigma' = \frac{-e^{-x}}{(1 + e^{-x})^2} (-1) \\ & & \Rightarrow & \sigma' = \frac{e^{-x}}{(1 + e^{-x})^2} \\ & & \Rightarrow & \sigma' = \frac{1}{1 + e^{-x}} \left(\frac{e^{-x}}{1 + e^{-x}} + \frac{1}{1 + e^{-x}} - \frac{1}{1 + e^{-x}} \right) \\ & & \Rightarrow & \sigma' = \frac{1}{1 + e^{-x}} \left(\frac{1 + e^{-x}}{1 + e^{-x}} - \frac{1}{1 + e^{-x}} \right) \\ & & \Rightarrow & \sigma' = \sigma(1 - \sigma) \quad \blacksquare & (\sigma(x) &= \frac{1}{1 + e^{-x}}) \end{aligned}$$

La fonction tanh

La fonction tangente hyperbolique (tanh) est une fonction d'activation dont le résultat est situé dans l'intervalle $[-1, 1]$:

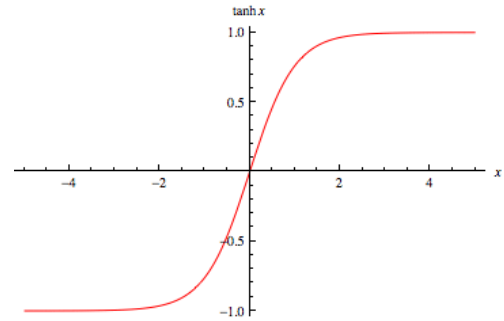
$$\text{Fonction : } \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$\text{Dérivée : } \frac{d \tanh(z)}{dz} = 1 - \tanh^2(z)$$

La fonction tanh est utilisé comme fonction d'activation à l'intérieur du réseau (couche cachée) par ce qu'elle admet des valeurs positives et négatives ce qui augmente la capacité d'apprentissage du réseau. Cependant, on l'évite à la dernière couche si le résultat est le fruit d'une classification (valeur de 0 ou 1).

Preuve :

En construction ...



<http://mathworld.wolfram.com/HyperbolicTangent.html>

Illustration des valeurs de la fonction tanh.

La fonction ReLU

La fonction Unité Linéaire Rectifiée

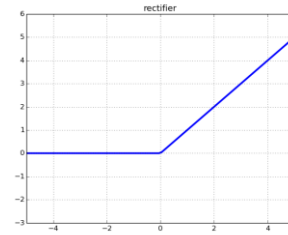
ReLU (rectified linear unit)

En construction ...

$$\text{Fonction : } \text{ReLU}(z) = \begin{cases} 0 & \text{si } z < 0 \\ z & \text{si } z \geq 0 \end{cases}$$

$$\text{ou bien } \text{ReLU}(z) = \max(0, z)$$

$$\text{Dérivée : } \frac{d\text{ReLU}(z)}{dz} = \begin{cases} 0 & \text{si } z < 0 \\ 1 & \text{si } z \geq 0 \end{cases}$$



<https://datascience.stackexchange.com/questions/5706/what-is-the-dying-relu-problem-in-neural-networks>

Illustration des valeurs de la fonction ReLU.

REMARQUE : (basée sur mon expérience personnelle)

- Puisque la méthode RELU fait bouger de beaucoup les poids et biais des couches, cette fonction d'activation nécessite un gradient avec plusieurs vecteurs (*mini-batch*). Autrement, il n'y aura pas de direction adéquate de prise lors de la descente du gradient et aucune bonne calibration des neurones seront prises.
- L'usage de la méthode ReLU nécessite la méthode SOFTMAX en dernière couche (fonction de classification), car les agrégations générées seront trop élevées et impossible à ajuster avec une fonction d'activation de classification comme SIGMOÏDE, car sa pente est trop faible à haute valeur d'agrégation.
- La fonction ReLU est réputé pour bien propager l'erreur en raison de sa pente de 1. Cependant, si sont activation est toujours négative peu importe le vecteur d'entrée, le neurone sera alors « désactivé » ou « mort ». Ainsi, il y aura moins de neurone pour mettre à jour le réseau. Ainsi, il est souhaitable d'utiliser un grand nombre de neurones dans une couche avec la fonction d'activation ReLU. Pour éviter ce problème, il existe des stratégies de normalisation comme la fonction *batch-normalisation* (voir section 6.3).

La fonction Leaky ReLU

$$\text{LReLU}(z) = \begin{cases} \varepsilon z & \text{si } z < 0 \\ z & \text{si } z \geq 0 \end{cases} \quad \text{ou bien} \quad \text{LReLU}(z) = \max(\varepsilon z, z)$$

avec $\varepsilon \in [0, 1[$, mais traditionnellement $\varepsilon = 0.1$

$$\frac{d\text{LReLU}(z)}{dz} = \begin{cases} \varepsilon & \text{si } z < 0 \\ 1 & \text{si } z \geq 0 \end{cases}$$

La fonction Randomized Leaky ReLU (2015)

On peut également faire une variante de cette fonction où ε peut être déterminé aléatoirement dans l'intervalle $\varepsilon \in [0, \varepsilon_{\max}]$ où ε_{\max} est inférieur à 1 pour chacun des neurones à activer par la fonction.

La fonction ELU

La fonction unité linéaire exponentiel

Exponential Linear Unit ELU

$$\text{ELU}(z) = \begin{cases} z & \text{si } z \geq 0 \\ \varepsilon(e^z - 1) & \text{si } z < 0 \end{cases}$$

avec $\varepsilon \in [0, 1[$, mais traditionnellement $\varepsilon = 0.1$

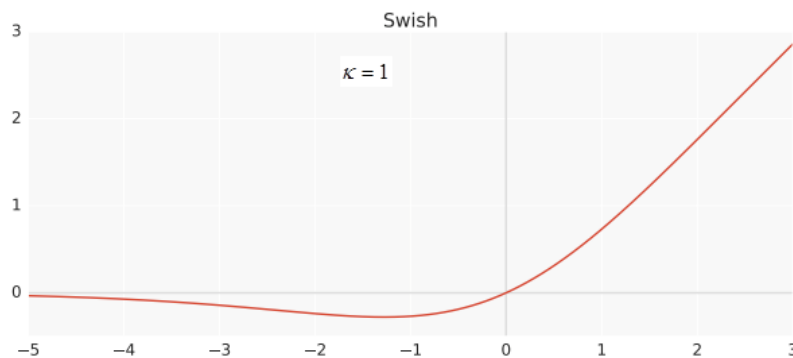
$$\frac{d\text{ELU}(z)}{dz} = \begin{cases} 1 & \text{si } z \geq 0 \\ \text{ELU}(z) + \varepsilon & \text{si } z < 0 \end{cases}$$

La fonction swish

$$\text{swish}(z, \kappa) \equiv \psi(z, \kappa) = \frac{z}{1 + e^{-\kappa z}}$$

$$\frac{\partial \psi}{\partial z} = \sigma + \kappa z \sigma (1 - \sigma) \quad \text{avec} \quad \sigma = \frac{1}{1 + e^{-\kappa z}}$$

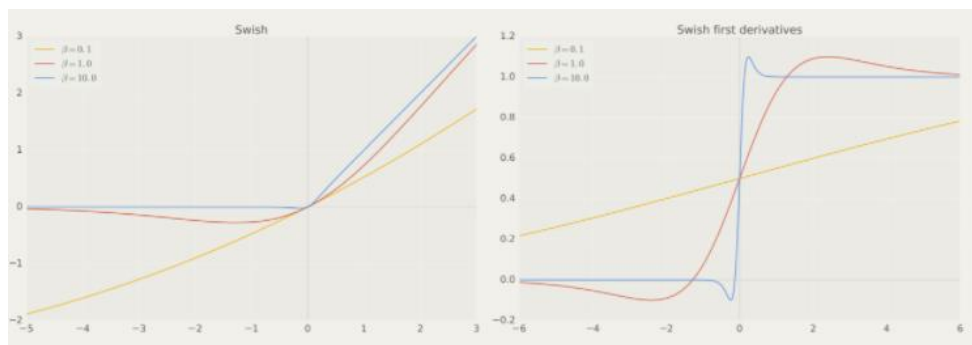
$$\text{et} \quad \frac{\partial \psi}{\partial \kappa} = z^2 \sigma (1 - \sigma)$$



<https://medium.com/@neuralnets/swish-activation-function-by-google-53e1ea86f820>

Le nom « swich » pour *sigmoid and wish* (développée par Google) permet à cette fonction de changer de comportement selon la valeur du paramètre κ qui peut être considéré par un réseau comme une constante ou comme un paramètre pouvant être entraîné (il y aura un paramètre κ pour chaque neurone u lors de l'activation d'une couche de neurones). Ainsi, entraîner le paramètre κ permet au réseau selon le type de donnée de choisir le format idéal de courbure de la fonction pour l'activation de l'agrégation.

- Si $\kappa = 1$, la fonction est équivalente à la fonction $f(z) = z \sigma(z)$ (*sigmoid-weighted Linear Unit*, SiL).
- Si $\kappa = 0$, la fonction est équivalent à la fonction linéaire $f(z) = \frac{z}{2}$.
- Si $\kappa \rightarrow \infty$, la fonction est équivalent à la fonction ReLU $f(z) = \max(0, z)$.



<https://sefiks.com/2018/08/21/swish-as-neural-networks-activation-function/>

Preuve : (version où $\beta \equiv \kappa$)

$$\begin{aligned} \frac{\partial \psi_u}{\partial z_u} &= \frac{\partial}{\partial z_u} \left(\frac{z_u}{1 + e^{-\beta z_u}} \right) \\ \Rightarrow \frac{\partial \psi_u}{\partial z_u} &= \frac{1}{1 + e^{-\beta z_u}} \frac{\partial z_u}{\partial z_u} + z_u \frac{\partial}{\partial z_u} \left(\frac{1}{1 + e^{-\beta z_u}} \right) \\ \Rightarrow \frac{\partial \psi_u}{\partial z_u} &= \frac{1}{1 + e^{-\beta z_u}} + z_u \frac{\partial}{\partial z_u} (1 + e^{-\beta z_u})^{-1} \\ \Rightarrow \frac{\partial \psi_u}{\partial z_u} &= \frac{1}{1 + e^{-\beta z_u}} + z_u (-1)(1 + e^{-\beta z_u})^{-2} \frac{\partial}{\partial z_u} (1 + e^{-\beta z_u}) \\ \Rightarrow \frac{\partial \psi_u}{\partial z_u} &= \frac{1}{1 + e^{-\beta z_u}} - \frac{z_u}{(1 + e^{-\beta z_u})^2} \frac{\partial}{\partial z_u} (e^{-\beta z_u}) \\ \Rightarrow \frac{\partial \psi_u}{\partial z_u} &= \frac{1}{1 + e^{-\beta z_u}} - \frac{z_u}{(1 + e^{-\beta z_u})^2} (e^{-\beta z_u}) \frac{\partial}{\partial z_u} (-\beta z_u) \\ \Rightarrow \frac{\partial \psi_u}{\partial z_u} &= \frac{1}{1 + e^{-\beta z_u}} - \frac{z_u}{(1 + e^{-\beta z_u})^2} (e^{-\beta z_u}) (-\beta) \frac{\partial z_u}{\partial z_u} \\ \Rightarrow \frac{\partial \psi_u}{\partial z_u} &= \frac{1}{1 + e^{-\beta z_u}} + \beta z_u \frac{e^{-\beta z_u}}{(1 + e^{-\beta z_u})^2} \end{aligned}$$

Si l'on effectue cette restructuration

$$\frac{e^{-\beta z_u}}{(1 + e^{-\beta z_u})^2} = \frac{e^{-\beta z_u} + 1 - 1}{(1 + e^{-\beta z_u})^2} = \frac{1}{1 + e^{-\beta z_u}} \frac{1 + e^{-\beta z_u} - 1}{1 + e^{-\beta z_u}} = \frac{1}{1 + e^{-\beta z_u}} \left(1 - \frac{1}{1 + e^{-\beta z_u}} \right)$$

et que l'on utilise l'expression

$$\sigma_u = \frac{1}{1 + e^{-\beta z_u}} ,$$

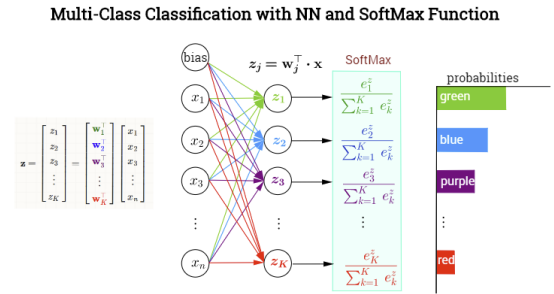
nous obtenons

$$\frac{\partial \psi_u}{\partial z_u} = \sigma_u + \beta z_u \sigma_u (1 - \sigma_u) \quad \blacksquare$$

$$\begin{aligned} \frac{\partial \psi_u}{\partial \beta} &= \frac{\partial}{\partial \beta} \left(\frac{z_u}{1 + e^{-\beta z_u}} \right) \\ \Rightarrow \frac{\partial \psi_u}{\partial \beta} &= z_u \frac{\partial}{\partial \beta} \left((1 + e^{-\beta z_u})^{-1} \right) \\ \Rightarrow \frac{\partial \psi_u}{\partial \beta} &= z_u (-1) (1 + e^{-\beta z_u})^{-2} \frac{\partial}{\partial \beta} (1 + e^{-\beta z_u}) \\ \Rightarrow \frac{\partial \psi_u}{\partial \beta} &= - \frac{z_u}{(1 + e^{-\beta z_u})^2} (e^{-\beta z_u}) \frac{\partial}{\partial \beta} (-\beta z_u) \\ \Rightarrow \frac{\partial \psi_u}{\partial \beta} &= - \frac{z_u e^{-\beta z_u}}{(1 + e^{-\beta z_u})^2} (-z_u) \\ \Rightarrow \frac{\partial \psi_u}{\partial \beta} &= z_u^2 \sigma_u (1 - \sigma_u) \end{aligned} \quad \frac{e^{-\beta z_u}}{(1 + e^{-\beta z_u})^2} = \sigma_u (1 - \sigma_u) \text{ et } \sigma_u = \frac{1}{1 + e^{-\beta z_u}}$$

La fonction softmax

La fonction d'activation *softmax* (fonction exponentielle normalisée) permet de générer une activation strictement positive dont la somme des activations sur l'ensemble des paramètres de la fonction sera égale à 1. Ainsi, la fonction *softmax* aura la forme d'une probabilité où l'activation de la fonction pour chacun de ses paramètres sera une probabilité parmi l'ensemble des probabilités. Cette fonction est régulièrement utilisée comme fonction d'activation à la toute dernière couche d'un réseau (fonction de classification) afin de réaliser une classification parmi plusieurs choix.



<http://rinterested.github.io/statistics/softmax.html>

Illustration de la distribution des probabilités générées par les activations des agrégations de la couche.

Avec sa nature probabiliste, l'activation correspondra à la probabilité que le réseau accorde à un choix de la classification où l'ensemble des choix possèdent une probabilité de 1.

Fonction	Dérivée partielle
$S_{\mu}(\vec{z}) = \frac{e^{z_{\mu}}}{\sum_{i=1}^N e^{z_i}}$ <p>avec $\vec{z} = (z_1, z_2, \dots, z_{\mu}, \dots, z_N)$</p> <p>et $\mu \in [1, N]$</p>	$\frac{\partial S_u(\vec{z})}{\partial z_k} = S_u(\delta_{uk} - S_k)$ <p>et $\delta_{kk} = 1, \delta_{uk} = 0$ (delta de Kronecker)</p>

Preuve :

Effectuons la différentielle de la fonction *softmax* $S_{\mu}(\vec{x})$ par rapport à la variable $x_{\mu} \in \vec{x}$:

$$dS_u(\vec{x}) = d \left(\frac{e^{x_u}}{\sum_{i=1}^N e^{x_i}} \right)$$

$$\Rightarrow dS_u(\vec{x}) = \sum_{k=1}^N \frac{\partial}{\partial x_k} \left(\frac{e^{x_u}}{\sum_{i=1}^N e^{x_i}} \right) dx_k$$

$$\Rightarrow dS_u(\vec{x}) = \frac{\partial}{\partial x_u} \left(\frac{e^{x_u}}{\sum_{i=1}^N e^{x_i}} \right) dx_u + \sum_{\substack{k=1 \\ k \neq u}}^N \frac{\partial}{\partial x_k} \left(\frac{e^{x_u}}{\sum_{i=1}^N e^{x_i}} \right) dx_k$$

$$\Rightarrow dS_u(\vec{x}) = \frac{1}{\left(\sum_{i=1}^N e^{x_i} \right)^2} \left[\sum_{i=1}^N e^{x_i} \frac{d}{dx_u} (e^{x_u}) - e^{x_u} \frac{\partial}{\partial x_u} \left(\sum_{i=1}^N e^{x_i} \right) \right] dx_u + \sum_{\substack{k=1 \\ k \neq u}}^N - \frac{e^{x_u}}{\left(\sum_{i=1}^N e^{x_i} \right)^2} \frac{\partial}{\partial x_k} \left(\sum_{i=1}^N e^{x_i} \right) dx_k$$

$$\Rightarrow dS_u(\bar{x}) = \frac{1}{\left(\sum_{i=1}^N e^{x_i}\right)^2} \left[\sum_{i=1}^N e^{x_i} e^{x_u} - e^{x_u} e^{x_u} \right] dx_u - \sum_{\substack{k=1 \\ k \neq u}}^N \frac{e^{x_u}}{\left(\sum_{i=1}^N e^{x_i}\right)^2} e^{x_k} dx_k$$

$$\Rightarrow dS_u(\bar{x}) = \frac{e^{x_u}}{\left(\sum_{i=1}^N e^{x_i}\right)^2} \left[\sum_{i=1}^N e^{x_i} - e^{x_u} \right] dx_u - \sum_{\substack{k=1 \\ k \neq u}}^N \frac{e^{x_u} e^{x_k}}{\left(\sum_{i=1}^N e^{x_i}\right)^2} dx_k$$

$$\Rightarrow dS_u(\bar{x}) = \frac{e^{x_u}}{\sum_{i=1}^N e^{x_i}} \left[1 - \frac{e^{x_u}}{\sum_{i=1}^N e^{x_i}} \right] dx_u - \sum_{\substack{k=1 \\ k \neq u}}^N \frac{e^{x_u}}{\sum_{i=1}^N e^{x_i}} \frac{e^{x_k}}{\sum_{i=1}^N e^{x_i}} dx_k$$

$$\Rightarrow dS_u(\bar{x}) = S_u (1 - S_u) dx_u - \sum_{\substack{k=1 \\ k \neq u}}^N S_u S_k dx_k$$

$$\Rightarrow dS_u(\bar{x}) = (S_u - S_u S_u) dx_u - \sum_{\substack{k=1 \\ k \neq u}}^N S_u S_k dx_k$$

La stratégie sera d'unifier nos deux dérivés par l'ajout d'un « rien stratégique ». Introduisons alors le delta de Kronecker⁴

$$\delta_{uk} \quad \text{où} \quad \delta_{kk} = 1 \text{ et } \delta_{uk} = 0 \quad \forall u, k \in [1, N] .$$

Ainsi, nous obtenons :

$$dS_u(\bar{x}) = (S_u - S_u S_u) dx_u - \sum_{\substack{k=1 \\ k \neq u}}^N S_u S_k dx_k \quad (\text{équation précédent})$$

$$\Rightarrow dS_u(\bar{x}) = (\delta_{uu} S_u - S_u S_u) dx_u + \sum_{\substack{k=1 \\ k \neq u}}^N (\delta_{uk} S_u - S_u S_k) dx_k$$

$$\Rightarrow dS_u(\bar{x}) = \sum_{k=1}^N (\delta_{uk} S_u - S_u S_k) dx_k$$

$$\Rightarrow dS_u(\bar{x}) = S_u \sum_{k=1}^N (\delta_{uk} - S_k) dx_k$$

L'expression de la différentielle prendra la forme

$$\frac{\partial S_u(\bar{x})}{\partial x_k} = S_u (\delta_{uk} - S_k)$$

sous la notation de la dérivée partielle. ■

⁴ Référence : https://en.wikipedia.org/wiki/Kronecker_delta

La fonction softmax numériquement stable

En raison de l'expression e^{z_μ} requise dans le calcul de la fonction *softmax*, une erreur numérique peut se propager lors d'une implémentation dans un programme informatique. Par exemple, lorsque $z_\mu = 40000$, la fonction *softmax* $S_\mu(\vec{z})$ peut engendrer un calcul numérique correspondant à ∞/∞ ce qui donne une valeur indéterminée (NaN, *not a number* en Java).

Afin d'augmenter la stabilité numérique de la fonction *softmax* $S_\mu(\vec{z})$, il est préférable de définir la fonction sous la forme suivante puisque la fonction *softmax* est invariante sous une translation m sur l'ensemble des composante z_μ du vecteur \vec{z} :

Fonction	Dérivée partielle
$S_\mu(\vec{z}) = \frac{e^{z_\mu - m}}{\sum_{i=1}^N e^{z_i - m}}$ <p>avec $\vec{z} = (z_1, z_2, \dots, z_\mu, \dots, z_N)$</p> <p>où $\mu \in [1, N]$</p> <p>et $m = \max(z_1, z_2, \dots, z_\mu, \dots, z_N)$</p>	$\frac{\partial S_u(\vec{z})}{\partial z_k} = S_u(\delta_{uk} - S_k)$ <p>et $\delta_{kk} = 1, \delta_{uk} = 0$ (delta de Kronecker)</p>

En réalisant une translation sur la valeur $m = \max(z_1, z_2, \dots, z_\mu, \dots, z_N)$ correspondant à la composante la plus grande du vecteur \vec{z} , nous obtenons la plus grande valeur numérique $e^0 = 1$ et la plus petite valeur numérique possible étant $e^{-\infty} = 0$. Ainsi, $S_\mu(\vec{z})$ ne pourra jamais donner une valeur numérique indéterminée.

Preuve :

Soit la translation c appliquée à la fonction *softmax*, démontrons que le résultat sera équivalent au résultat sans translation :

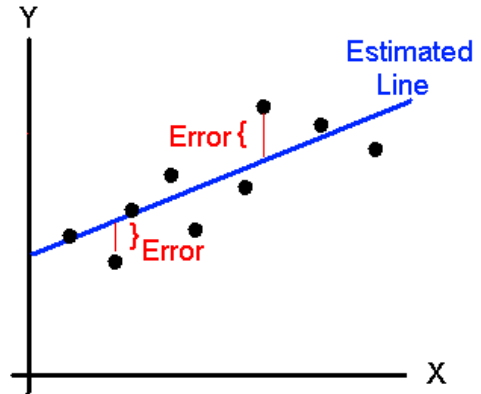
$$S_\mu(\vec{z}) = \frac{e^{z_\mu}}{\sum_{i=1}^N e^{z_i}} \Rightarrow S_\mu(\vec{z}) = \frac{e^c}{e^c} \frac{e^{z_\mu}}{\sum_{i=1}^N e^{z_i}} \quad (\text{Multiplier par } \frac{e^c}{e^c})$$

$$\Rightarrow S_\mu(\vec{z}) = \frac{e^{z_\mu + c}}{\sum_{i=1}^N e^{z_i + c}} \quad (\text{Propriété des exponentiels : } e^a e^b = e^{a+b})$$

$$\Rightarrow S_\mu(\vec{z}) = S_\mu(\vec{z} + \vec{c}) \quad \vec{c} = (c, c, c, \dots, c) \text{ où } c_u = c$$

La fonction d'erreur

La fonction d'erreur⁵ C_i évalue la différence entre les valeurs calculées $a^{(L)}$ par le réseau de neurones à partir d'un vecteur d'entrée $a^{(*)}$ avec les valeurs finales attendues y_i associées au vecteur d'entrée. Si le réseau reproduit exactement les valeurs attendues pour un vecteur d'entrée $a^{(*)}$, alors l'erreur sera nulle et aucun apprentissage ne sera possible pour ce vecteur. Autrement, la fonction d'erreur sera utilisée pour mettre à jour les poids $\tilde{w}_{vu}^{(k)}$ et biais $b_u^{(k)}$ du réseau.



Mathématiquement, la fonction d'erreur prendra la forme de l'équation suivante :

$$C = C(a_u^{(L)}, y_u) \text{ et } u \in [0, N^{(L)} - 1]$$

http://wiki.fast.ai/images/5/55/Linear_line_w_cost_function.png

L'objectif de la fonction d'erreur est d'évaluer l'écart entre une valeur calculée et une valeur désirée (*Estimated Line*).

Il existe plusieurs fonctions d'erreur. Le choix dépendra de la fonction d'activation utilisée à la dernière couche du réseau.

La différentielle de la fonction d'erreur

L'expression de la différentielle de la fonction d'erreur

$$C = C(a_u^{(L)}, y_u) \quad \text{où} \quad u \in [0 \dots N^{(L)} - 1]$$

est égale à l'expression suivante :

$$dC = \sum_{u=0}^{N^{(L)}-1} \frac{\partial C}{\partial a_u^{(L)}} da_u^{(L)}$$

(Différentielle de la fonction d'erreur)

Preuve :

Soit la fonction d'erreur

$$C = C(a_i^{(L)}, y_i)$$

avec $i \in [0 \dots N^{(L)} - 1]$ étant l'indice de sortie de l'activation. Appliquons la différentielle sur la fonction d'erreur ce qui nous donnera

$$dC = \sum_{i=0}^{N^{(L)}-1} \left(\frac{\partial C}{\partial a_i^{(L)}} da_i^{(L)} + \frac{\partial C}{\partial y_i} dy_i \right) .$$

Puisque $dy_i = 0 \quad \forall i \in [0 \dots N^{(L)} - 1]$, nous obtenons

$$dC = \sum_{i=0}^{N^{(L)}-1} \frac{\partial C}{\partial a_i^{(L)}} da_i^{(L)} \quad . \blacksquare$$

⁵ Il existe plusieurs synonymes à cette fonction : *Error*, *Cost*, *Loss*.

La fonction d'erreur linéaire

En construction ...

Fonction	Dérivée partielle
$C = \sum_{i=0}^{N^{(L)}-1} C_i$ avec $C_i = a_i^{(L)} - y_i$	$\frac{\partial C}{\partial a_i^{(L)}} = 1$

où y_i : Valeur attendue pour le neurone i .

$a_i^{(L)}$: Valeur obtenue par le réseau pour le neurone i .

La fonction d'erreur quadratique

En construction ...

Cette fonction porte également le nom populaire « d'énergie » pour son comportement quadratique tel que l'énergie cinétique $K = 0,5 mv^2$.

Fonction	Dérivée partielle
$C = \sum_{i=0}^{N^{(L)}-1} C_i$ avec $C_i = \frac{1}{2} (a_i^{(L)} - y_i)^2$	$\frac{\partial C}{\partial a_i^{(L)}} = a_i^{(L)} - y_i$

où y_i : Valeur attendue pour le neurone i .

$a_i^{(L)}$: Valeur obtenue par le réseau pour le neurone i .

La fonction d'erreur *cross entropy*

En construction ...

Cette fonction est utilisée lorsque y_i et $a_i^{(L)}$ correspondent à des distributions de probabilité. On utilise habituellement cette fonction d'erreur lorsque la dernière activation est la fonction *softmax*. Ce choix permet des optimisations dans les calculs de l'apprentissage.

Fonction	Dérivée partielle
$C = \sum_{i=0}^{N^{(L)}-1} C_i$ avec $C_i = -y_i \ln(a_i^{(L)})$	$\frac{\partial C}{\partial a_i^{(L)}} = -\frac{y_i}{a_i^{(L)}}$

où y_i : Valeur probabiliste attendue pour le neurone i .

$a_i^{(L)}$: Valeur probabiliste obtenue par le réseau pour le neurone i .

L'apprentissage

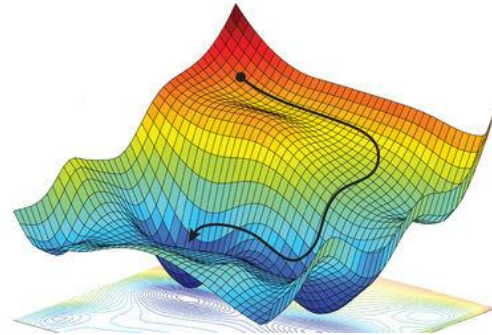
Le gradient de la fonction d'erreur

Le gradient de la fonction d'erreur ∇C qui dépend des dérivées partielles

$$\frac{\partial C}{\partial w_{uv}^{(k)}} \text{ et } \frac{\partial C}{\partial b_u^{(k)}}$$

par rapport au poids $w_{uv}^{(k)}$ et biais $b_u^{(k)}$ d'un neurone u d'une couche k permet d'évaluer à quel point ces paramètres du réseau influencent le calcul de l'activation du réseau pour un vecteur d'entrée $a^{(*)}$. Le calcul du gradient dépendra du type de réseau, du choix de la fonction d'erreur et de l'activation du vecteur d'entrée choisi. Sa valeur changera après chaque modification du réseau (modification de $w_{uv}^{(k)}$ et $b_u^{(k)}$) et tendra vers zéro lorsque la fonction d'erreur sera minimisée.

Puisque l'objectif de l'apprentissage sera de réduire la valeur de la fonction d'erreur pour l'ensemble des données d'entraînement, évaluer le gradient permettra de savoir la direction et le dosage du changement à apporter aux paramètres $w_{uv}^{(k)}$ et $b_u^{(k)}$ afin de réduire l'erreur. Puisque le gradient pointe vers une hausse de la fonction d'erreur, il faudra modifier les paramètres $w_{uv}^{(k)}$ et $b_u^{(k)}$ dans le sens contraire du gradient.



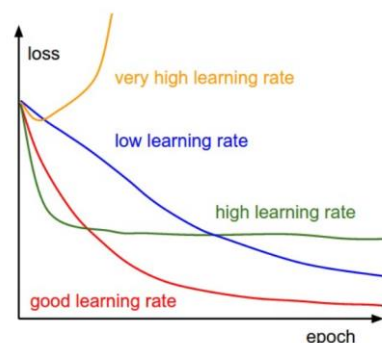
<https://www.sciencemag.org/news/2018/05/ai-researchers-allege-machine-learning-alchemy>

Le gradient de la fonction d'erreur permet de savoir dans quelle direction il faut modifier les paramètres du réseau afin de réduire l'erreur sur l'ensemble des données d'entraînement.

Le déplacement dans la fonction d'erreur sera orienté par le gradient de la fonction d'erreur ∇C , mais l'ampleur du déplacement sera pondéré par un facteur portant le nom de taux d'apprentissage (*learning rate*).

Le taux d'apprentissage α influence fortement le rythme de l'apprentissage. Cependant, il influence également la qualité de la convergence de la solution. Il est traditionnellement difficile à déterminer un taux d'apprentissage idéal. Les taux traditionnellement utilisés sont les suivants :

0.001 , 0.003 , 0.01 , 0.03 , 0.1 et 0.3



<https://towardsdatascience.com/gradient-descent-algorithm-and-its-variants-10f652806a3>

Graphique de l'évolution de la précision du réseau en fonction du nombre de donnée analysé dans l'apprentissage.

L'algorithme de l'apprentissage

Afin que le réseau puisse réaliser un « apprentissage », il faut modifier l'ensemble de ses poids $w_{uv}^{(k)}$ et biais $b_u^{(k)}$ judicieusement afin de minimiser la fonction d'erreur C pour un ensemble de données d'entraînement. Si ces données sont préalablement connues⁶, le réseau réalise un « apprentissage supervisé ». Cette méthode itérative porte le nom de rétropropagation du gradient (*backpropagation*). et $b_u^{(k)}$

À partir du gradient de la fonction d'erreur ∇C , nous pouvons modifier les paramètres $w_{uv}^{(k)}$ et $b_u^{(k)}$ du réseau grâce aux équations suivantes :

Au réseau W	Poids $\tilde{w}_{vu}^{(k)}$	Biais $b_u^{(k)}$
$W^{\text{mod}} = W - \alpha \nabla C$	$\tilde{w}_{vu}^{\text{mod}(k)} = \tilde{w}_{vu}^{(k)} - \alpha \frac{\partial C}{\partial \tilde{w}_{vu}^{(k)}}$	$b_u^{\text{mod}(k)} = b_u^{(k)} - \alpha \frac{\partial C}{\partial b_u^{(k)}}$

où

$\tilde{w}_{vu}^{\text{mod}(k)}, b_u^{\text{mod}(k)}$: Valeur modifiée d'un poids/biais (après apprentissage).

$\tilde{w}_{vu}^{(k)}, b_u^{(k)}$: Valeur d'un poids/biais à l'activation (avant apprentissage).

$\frac{\partial C}{\partial \tilde{w}_{vu}^{(k)}}, \frac{\partial C}{\partial b_u^{(k)}}$: Gradient de la fonction d'erreur pour un poids/biais de couche k .

α : Le taux d'apprentissage (*learning rate*) $\alpha > 0$.

⁶ Si les données d'entraînement ne sont pas connues, l'entraînement sera « non supervisé ». L'apprentissage par renforcement propose une technique pour générer des données d'entraînement dont la qualité augmente avec la progression de l'entraînement.

La propagation du gradient

Le gradient de la fonction d'erreur C noté ∇C par rapport aux paramètres de poids $\tilde{w}_{vu}^{(k)}$ et biais $b_u^{(k)}$ de toutes les couches du réseau dépend du choix de la fonction d'erreur, mais plus explicitement du choix et de l'ordre des fonctions à l'intérieur du réseau.

Dans un réseau simple⁷, il n'y aura que des fonctions d'agrégation

$$z_u^{(k)} = z_u^{(k)}(a_v^{(k-1)}, \tilde{w}_{vu}^{(k)}, b_u^{(k)})$$

et des fonctions d'activation à un neurone⁸

$$a_u^{(k)} = a_u^{(k)}(z_u^{(k)})$$

dans l'ordre

$$a^{(*)} \rightarrow z^{(0)} \rightarrow a^{(0)} \rightarrow z^{(1)} \rightarrow a^{(1)} \rightarrow \dots \rightarrow z^{(k)} \rightarrow a^{(k)} \rightarrow z^{(k+1)} \rightarrow a^{(k+1)} \rightarrow \dots \rightarrow z^{(L)} \rightarrow a^{(L)}$$

où

$a^{(*)}$ sont les entrées du réseau et $a^{(L)}$ sont les sorties du réseau (l'activation du réseau).

En exploitant le concept mathématique de la différentielle d'une fonction et de la dérivée en chaîne

$$\frac{dy}{dx} = \frac{dy}{du} \cdot \frac{du}{dx},$$

nous obtenons l'expression suivante pour la différentielle de la fonction d'erreur :

$$dC = \sum_{i=0}^{N^{(L)}-1} \frac{\partial C}{\partial a_i^{(L)}} \frac{\partial a_i^{(L)}}{\partial z_i^{(L)}} \left(\sum_{j=0}^{N^{(L-1)}-1} \frac{\partial z_i^{(L)}}{\partial \tilde{w}_{ji}^{(L)}} d\tilde{w}_{ji}^{(L)} + \frac{\partial z_i^{(L)}}{\partial b_i^{(L)}} db_i^{(L)} \right) + \sum_{j=0}^{N^{(L-1)}-1} \frac{\partial z_i^{(L)}}{\partial a_j^{(L-1)}} \frac{\partial a_j^{(L-1)}}{\partial z_j^{(L-1)}} \left(\sum_{k=0}^{N^{(L-2)}-1} \frac{\partial z_j^{(L-1)}}{\partial \tilde{w}_{kj}^{(L-1)}} d\tilde{w}_{kj}^{(L-1)} + \frac{\partial z_j^{(L-1)}}{\partial b_j^{(L-1)}} db_j^{(L-1)} \right) + \sum_{k=0}^{N^{(L-2)}-1} \frac{\partial z_j^{(L-1)}}{\partial a_k^{(L-2)}} da_k^{(L-2)} \right)$$

(Propagation du gradient de profondeur $n = 1$)

Cette lourde équation permet d'accéder aux équations nécessaires pour évaluer l'ensemble des termes

$$\frac{\partial C}{\partial \tilde{w}_{vu}^{(k)}} \quad \text{et} \quad \frac{\partial C}{\partial b_u^{(k)}}$$

permettant de connaître comment on doit modifier tous les paramètres $\tilde{w}_{vu}^{(k)}$ et $b_u^{(k)}$ à l'aide de l'équation de l'apprentissage

$$\tilde{w}_{vu}^{\text{mod}(k)} = \tilde{w}_{vu}^{(k)} - \alpha \frac{\partial C}{\partial \tilde{w}_{vu}^{(k)}} \quad \text{et} \quad b_u^{\text{mod}(k)} = b_u^{(k)} - \alpha \frac{\partial C}{\partial b_u^{(k)}}$$

du réseau pour que celui-ci converge vers la fonction que l'on désire reproduire. Pour bien suivre la séquence des modifications, nous procéderons par plusieurs étapes qui seront présentées dans la suite de ce texte.

⁷ Un réseau avec normalisation aura une fonction supplémentaire à traiter dans le calcul de la dérivée partielle de la fonction d'erreur.

⁸ Une fonction d'activation à un neurone nécessite seulement une valeur au calcul comme la fonction sigmoïde.

Preuve :

Soit la fonction d'erreur $C = C(a_i^{(L)}, y_i)$, nous obtenons alors la différentielle suivante :

$$dC = \sum_{i=0}^{N^{(L)}-1} \frac{\partial C}{\partial a_i^{(L)}} da_i^{(L)}$$

Avec une fonction d'activation $a_i^{(L)} = a_i^{(L)}(z_i^{(L)})$ à un neurone telle que la différentielle est égale à

$$da_i^{(L)} = \frac{\partial a_i^{(L)}}{\partial z_i^{(L)}} dz_i^{(L)},$$

ceci qui donnera la différentielle suivante :

$$dC = \sum_{i=0}^{N^{(L)}-1} \frac{\partial C}{\partial a_i^{(L)}} \frac{\partial a_i^{(L)}}{\partial z_i^{(L)}} dz_i^{(L)}$$

Avec la fonction d'agrégation $z_i^{(L)} = z_i^{(L)}(a_j^{(L-1)}, \tilde{w}_{ji}^{(L)}, b_i^{(L)})$ telle que la différentielle est égale à

$$dz_i^{(L)} = \sum_{j=0}^{N^{(L-1)}-1} \left(\frac{\partial z_i^{(L)}}{\partial a_j^{(L-1)}} da_j^{(L-1)} + \frac{\partial z_i^{(L)}}{\partial \tilde{w}_{ji}^{(L)}} d\tilde{w}_{ji}^{(L)} \right) + \frac{\partial z_i^{(L)}}{\partial b_i^{(L)}} db_i^{(L)}$$

ceci nous donnera la différentielle

$$dC = \sum_{i=0}^{N^{(L)}-1} \frac{\partial C}{\partial a_i^{(L)}} \frac{\partial a_i^{(L)}}{\partial z_i^{(L)}} \left(\sum_{j=0}^{N^{(L-1)}-1} \left(\frac{\partial z_i^{(L)}}{\partial a_j^{(L-1)}} da_j^{(L-1)} + \frac{\partial z_i^{(L)}}{\partial \tilde{w}_{ji}^{(L)}} d\tilde{w}_{ji}^{(L)} \right) + \frac{\partial z_i^{(L)}}{\partial b_i^{(L)}} db_i^{(L)} \right)$$

que nous pouvons réécrire de la façon suivante :

$$dC = \sum_{i=0}^{N^{(L)}-1} \frac{\partial C}{\partial a_i^{(L)}} \frac{\partial a_i^{(L)}}{\partial z_i^{(L)}} \left(\sum_{j=0}^{N^{(L-1)}-1} \frac{\partial z_i^{(L)}}{\partial \tilde{w}_{ji}^{(L)}} d\tilde{w}_{ji}^{(L)} + \frac{\partial z_i^{(L)}}{\partial b_i^{(L)}} db_i^{(L)} \right) + \sum_{j=0}^{N^{(L-1)}-1} \frac{\partial z_i^{(L)}}{\partial a_j^{(L-1)}} da_j^{(L-1)}$$

(Propagation du gradient de niveau $n = 0$)

Si nous continuons à développer la différentielle au niveau $n = 1$, utilisons la fonction d'activation $a_j^{(L-1)} = a_j^{(L-1)}(z_j^{(L-1)})$ telle que la différentielle est égale à

$$da_j^{(L-1)} = \frac{\partial a_j^{(L-1)}}{\partial z_j^{(L-1)}} dz_j^{(L-1)}$$

et ceci nous donne l'expression suivante :

$$dC = \sum_{i=0}^{N^{(L)}-1} \frac{\partial C}{\partial a_i^{(L)}} \frac{\partial a_i^{(L)}}{\partial z_i^{(L)}} \left(\sum_{j=0}^{N^{(L-1)}-1} \frac{\partial z_i^{(L)}}{\partial \tilde{w}_{ji}^{(L)}} d\tilde{w}_{ji}^{(L)} + \frac{\partial z_i^{(L)}}{\partial b_i^{(L)}} db_i^{(L)} \right) + \sum_{j=0}^{N^{(L-1)}-1} \frac{\partial z_i^{(L)}}{\partial a_j^{(L-1)}} \frac{\partial a_j^{(L-1)}}{\partial z_j^{(L-1)}} dz_j^{(L-1)}$$

Avec la fonction d'agrégation $z_j^{(L-1)} = z_j^{(L-1)}(a_k^{(L-2)}, \tilde{w}_{kj}^{(L-1)}, b_j^{(L-1)})$ telle que la différentielle est égale à

$$dz_j^{(L-1)} = \sum_{k=0}^{N^{(L-2)}-1} \left(\frac{\partial z_j^{(L-1)}}{\partial a_k^{(L-2)}} da_k^{(L-2)} + \frac{\partial z_j^{(L-1)}}{\partial \tilde{w}_{kj}^{(L-1)}} d\tilde{w}_{kj}^{(L-1)} \right) + \frac{\partial z_j^{(L-1)}}{\partial b_j^{(L-1)}} db_j^{(L-1)}$$

ceci nous donnera la différentielle suivante après réorganisation des termes :

$$dC = \sum_{i=0}^{N^{(L)}-1} \frac{\partial C}{\partial a_i^{(L)}} \frac{\partial a_i^{(L)}}{\partial z_i^{(L)}} \left(\sum_{j=0}^{N^{(L-1)}-1} \frac{\partial z_i^{(L)}}{\partial \tilde{w}_{ji}^{(L)}} d\tilde{w}_{ji}^{(L)} + \frac{\partial z_i^{(L)}}{\partial b_i^{(L)}} db_i^{(L)} \right) + \sum_{j=0}^{N^{(L-1)}-1} \frac{\partial z_i^{(L)}}{\partial a_j^{(L-1)}} \frac{\partial a_j^{(L-1)}}{\partial z_j^{(L-1)}} \left(\sum_{k=0}^{N^{(L-2)}-1} \frac{\partial z_j^{(L-1)}}{\partial \tilde{w}_{kj}^{(L-1)}} d\tilde{w}_{kj}^{(L-1)} + \frac{\partial z_j^{(L-1)}}{\partial b_j^{(L-1)}} db_j^{(L-1)} \right) + \sum_{k=0}^{N^{(L-2)}-1} \frac{\partial z_j^{(L-1)}}{\partial a_k^{(L-2)}} da_k^{(L-2)}$$

(Propagation du gradient de niveau $n = 1$)

En fixant les indices i et j , nous pouvons évaluer les coefficients devant les termes $d\tilde{w}_{ji}^{(L)}$ et $db_i^{(L)}$ ce qui correspondra aux dérivées partielles de la fonction d'erreur

$$\frac{\partial C}{\partial \tilde{w}_{ji}^{(L)}} \quad \text{et} \quad \frac{\partial C}{\partial b_i^{(L)}}$$

par rapport aux paramètres de la couche L du réseau.

En effectuant la sommation sur l'indice i et en fixant l'indice j et k , nous pouvons évaluer les coefficients devant les termes $d\tilde{w}_{kj}^{(L-1)}$ et $db_j^{(L-1)}$ ce qui correspondra aux dérivées partielles de la fonction d'erreur

$$\frac{\partial C}{\partial \tilde{w}_{jk}^{(L-1)}} \quad \text{et} \quad \frac{\partial C}{\partial b_j^{(L-1)}}$$

par rapport aux paramètres de la couche $L-1$ du réseau. ■

Les équations de la propagation du gradient

L'équation de la propagation de l'erreur de la fonction d'erreur

Lorsque l'on évalue le gradient de la fonction d'erreur C , nous allons devoir calculer comment une erreur évaluée sur la fonction d'erreur sera propagée à l'ensemble des paramètres du réseau.

Ainsi, en appliquant la différentielle à notre fonction d'erreur C , nous pouvons évaluer l'expression suivante pour évaluer l'erreur $\Delta_u^{C(L)}$ propagée par la fonction d'erreur C situé à la couche L du réseau pour le neurone de sortie u du réseau :

À la couche L	
$\Delta_u^{C(L)} = \frac{\partial C}{\partial a_u^{(L)}} \text{ avec } u \in [0, N^{(L)} - 1]$	

Preuve :

Soit la fonction d'erreur $C = C(a_u^{(L)}, y_u)$ où y_u est la valeur attendue que le réseau devrait générer pour $a_u^{(L)}$ avec l'activation du vecteur d'entrée $a^{(*)}$. Appliquons la différentielle à cette fonction ce qui nous donnera le terme de propagation de l'erreur de la fonction d'erreur pour le neurone u :

$$\begin{aligned} dC = dC(a_u^{(L)}, y_u) &\Rightarrow dC = \sum_{u=0}^{N^{(L)}-1} \left(\frac{\partial C}{\partial a_u^{(L)}} da_u^{(L)} + \frac{\partial C}{\partial y_u} dy_u \right) \\ &\Rightarrow dC = \sum_{u=0}^{N^{(L)}-1} \frac{\partial C}{\partial a_u^{(L)}} da_u^{(L)} && (dy_u = 0, \text{ car la valeur attendue ne varie pas}) \\ &\Rightarrow dC = \sum_{u=0}^{N^{(L)}-1} \Delta_u^{C(L)} da_u^{(L)} \quad \blacksquare && (\Delta_u^{C(L)} = \frac{\partial C}{\partial a_u^{(L)}}) \end{aligned}$$

L'équation de la propagation de l'erreur de la fonction d'activation

Lorsque l'on évalue le gradient de la fonction d'erreur C , l'erreur propagée par la fonction d'erreur sera propagée à nouveau par la fonction d'activation, car celle-ci devra propager son erreur à la fonction d'agrégation qui elle aura une influence sur les poids $\tilde{w}_{vu}^{(k)}$ et biais $b_u^{(k)}$.

Ainsi, en appliquant la différentielle à notre fonction d'erreur C , nous pouvons évaluer l'expression suivante pour évaluer l'erreur $\Delta_u^{a(k)}$ propagée par la fonction d'activation $a^{(k)}$ situé à la couche k du réseau pour le neurone de sortie u de la couche k du réseau :

À la couche L	À la couche k
$\Delta_u^{a(L)} = \Delta_u^{C(L)} \frac{\partial a_u^{(k)}}{\partial z_u^{(k)}} \text{ avec } u \in [0, N^{(k)} - 1]$	$\Delta_u^{a(k)} = \Delta_u^{C(L)} \frac{\partial a_u^{(k)}}{\partial z_u^{(k)}} \text{ avec } u \in [0, N^{(k)} - 1]$

Preuve :

Soit la différentielle de la fonction d'erreur C dont l'erreur propagée à la fonction d'activation $a_u^{(k)}$ de niveau de propagation $k = L - n$ est égale à l'expression

$$dC = \sum_{u=0}^{N^{(k)}-1} \Delta_u^{(k)} da_u^{(k)} \quad \text{avec} \quad da_u^{(k)} = \frac{\partial a_u^{(k)}}{\partial z_u^{(k)}} dz_u^{(k)}$$

où $\Delta_u^{(k)}$ est la somme des erreurs propagée depuis la fonction d'erreur vers l'intérieur du réseau pour l'ensemble des fonctions entre la fonction d'erreur C et la fonction d'activation $a_u^{(k)}$.

Le paramètre n est le niveau de profondeur dans la propagation de l'erreur dans le réseau. Il est à noter que $n = 0$ lorsque l'erreur est propagée par la dernière couche du réseau et $n = L$ lorsque l'erreur est propagée à la 1^{re} couche du réseau.

Développons l'expression afin d'évaluer l'erreur qui sera propagée par la fonction $a_u^{(L)} = a_u^{(L)}(z_u^{(L)})$ pour chaque neurone u de cette fonction à la fonction d'agrégation $z_u^{(k)}$:

$$\begin{aligned} dC = \sum_{u=0}^{N^{(k)}-1} \Delta_u^{(k)} da_u^{(k)} &\Rightarrow dC = \sum_{u=0}^{N^{(k)}-1} \Delta_u^{(k)} \frac{\partial a_u^{(k)}}{\partial z_u^{(k)}} dz_u^{(k)} && (da_u^{(L)} = \frac{\partial a_u^{(L)}}{\partial z_u^{(L)}} dz_u^{(L)}) \\ &\Rightarrow dC = \sum_{u=0}^{N^{(k)}-1} \Delta_u^{(k)} dz_u^{(k)} \quad \blacksquare && (\Delta_u^{(k)} = \Delta_u^{(k)} \frac{\partial a_u^{(k)}}{\partial z_u^{(k)}}) \end{aligned}$$

L'équation de la propagation de l'erreur de la fonction d'activation à plusieurs neurones

Pour évaluer la propagation de l'erreur d'une fonction d'activation à plusieurs neurones $a_u^{full(k)} = f_u^{full(k)}(\vec{z}^{(k)})$, il ne suffit que d'effectuer une sommation sur l'ensemble des neurones sous la forme suivante puisque toutes les valeurs d'activation ont été influencées par toutes les agrégations $z_u^{(k)}$:

À la couche L	À la couche k
$\Delta_u^{C-full(k)} = \sum_{p=0}^{N^{(k)}-1} \Delta_p^{C(L)} \frac{\partial a_p^{(k)}}{\partial z_u^{(k)}}$ <p>avec $u, p \in [0, N^{(k)} - 1]$</p>	$\Delta_u^{a-full(k)} = \sum_{p=0}^{N^{(k)}-1} \Delta_p^{(k)} \frac{\partial a_p^{(k)}}{\partial z_u^{(k)}}$ <p>avec $u, p \in [0, N^{(k)} - 1]$</p>

Preuve :

Soit la différentielle de la fonction d'erreur C dont l'erreur propagée à la fonction d'activation

$$a_u^{full(k)} = a_u^{full(k)}(z_0^{(k)}, z_1^{(k)}, \dots, z_p^{(k)}, \dots, z_{N^{(k)}-1}^{(k)})$$

de niveau de propagation $k = L - n$ est égale à l'expression

$$dC = \sum_{u=0}^{N^{(k)}-1} \Delta_u^{(k)} da_u^{full(k)} \quad \text{où} \quad da_u^{full(k)} = \sum_{p=0}^{N^{(k)}-1} \frac{\partial a_u^{full(k)}}{\partial z_p^{(k)}} dz_p^{(k)}$$

correspond à la différentielle de la fonction d'activation à plusieurs neurones.

Développons l'expression afin d'évaluer l'erreur qui sera propagée par la fonction $a_u^{full(k)}$ pour chaque neurone u de cette fonction à la fonction d'agrégation $z_u^{(k)}$:

$$\begin{aligned}
 dC &= \sum_{u=0}^{N^{(k)}-1} \Delta_u^{(k)} da_u^{(k)} &\Rightarrow & dC = \sum_{u=0}^{N^{(k)}-1} \Delta_u^{(k)} \sum_{p=0}^{N^{(k)}-1} \frac{\partial a_u^{(k)}}{\partial z_p^{(k)}} dz_p^{(k)} && (da_u^{full(k)} = \sum_{p=0}^{N^{(k)}-1} \frac{\partial a_u^{full(k)}}{\partial z_p^{(k)}} dz_p^{(k)}) \\
 & &\Rightarrow & dC = \sum_{p=0}^{N^{(k)}-1} \sum_{u=0}^{N^{(k)}-1} \Delta_u^{(k)} \frac{\partial a_u^{(k)}}{\partial z_p^{(k)}} dz_p^{(k)} && (\text{Factoriser la sommation sur } p) \\
 & &\Rightarrow & dC = \sum_{u=0}^{N^{(k)}-1} \sum_{p=0}^{N^{(k)}-1} \Delta_p^{(k)} \frac{\partial a_p^{(k)}}{\partial z_u^{(k)}} dz_u^{(k)} && (\text{Inverser indice } p \leftrightarrow u) \\
 & &\Rightarrow & dC = \sum_{u=0}^{N^{(k)}-1} \Delta_u^{a-full(k)} dz_u^{(k)} \blacksquare && (\Delta_u^{a-full(k)} = \sum_{p=0}^{N^{(k)}-1} \Delta_p^{(k)} \frac{\partial a_p^{(k)}}{\partial z_u^{(k)}})
 \end{aligned}$$

L'équation de la propagation de l'erreur de la fonction d'agrégation

Pour évaluer la propagation de l'erreur d'une fonction d'agrégation $z_u^{(k)}$ et le gradient de la fonction d'erreur par rapport au poids et biais d'une couche, nous pouvons évaluer les expressions suivantes :

Propagation de l'erreur	Gradient d'un poids	Gradient d'un biais
$\Delta_v^{z(k)} = \sum_{u=0}^{N^{(k)}-1} \Delta_u^{(k)} \frac{\partial z_u^{(k)}}{\partial a_v^{(k-1)}}$	$\frac{\partial C}{\partial \tilde{w}_{vu}^{(k)}} = \Delta_u^{(k)} \frac{\partial z_u^{(k)}}{\partial \tilde{w}_{vu}^{(k)}}$	$\frac{\partial C}{\partial b_u^{(k)}} = \Delta_u^{(k)} \frac{\partial z_u^{(k)}}{\partial b_u^{(k)}}$

$$\text{où } u \in [0, N^{(k)} - 1] \quad \text{et} \quad v \in [0, N^{(k-1)} - 1]$$

Preuve :

Soit la différentielle de la fonction d'erreur C dont l'erreur propagée à la fonction d'agrégation $z_u^{(k)}$ de niveau de propagation $k = L - n$ est égale à l'expression

$$dC = \sum_{u=0}^{N^{(k)}-1} \Delta_u^{(k)} dz_u^{(k)} \quad \text{où} \quad dz_u^{(k)} = \sum_{v=0}^{N^{(k-1)}-1} \left(\frac{\partial z_u^{(k)}}{\partial a_v^{(k-1)}} da_v^{(k-1)} + \frac{\partial z_u^{(k)}}{\partial \tilde{w}_{vu}^{(k)}} d\tilde{w}_{vu}^{(k)} \right) + \frac{\partial z_u^{(k)}}{\partial b_u^{(k)}} db_u^{(L)}$$

Développons l'expression afin d'évaluer l'erreur qui sera propagée par la fonction $z_u^{(k)}$ pour chaque neurone u de cette fonction à la fonction d'activation d'entrée $a_v^{(k-1)}$ ayant un nombre de neurones v différents :

$$\begin{aligned}
 dC &= \sum_{u=0}^{N^{(k)}-1} \Delta_u^{(k)} dz_u^{(k)} \\
 \Rightarrow dC &= \sum_{u=0}^{N^{(k)}-1} \Delta_u^{(k)} \left(\sum_{v=0}^{N^{(k-1)}-1} \left(\frac{\partial z_u^{(k)}}{\partial a_v^{(k-1)}} da_v^{(k-1)} + \frac{\partial z_u^{(k)}}{\partial \tilde{w}_{vu}^{(k)}} d\tilde{w}_{vu}^{(k)} \right) + \frac{\partial z_u^{(k)}}{\partial b_u^{(k)}} db_u^{(L)} \right)
 \end{aligned}$$

$$\Rightarrow dC = \sum_{u=0}^{N^{(k)}-1} \sum_{v=0}^{N^{(k-1)}-1} \Delta_u^{(k)} \frac{\partial z_u^{(k)}}{\partial \tilde{w}_{vu}^{(k)}} d\tilde{w}_{vu}^{(k)} + \sum_{u=0}^{N^{(k)}-1} \sum_{v=0}^{N^{(k-1)}-1} \Delta_u^{(k)} \frac{\partial z_u^{(k)}}{\partial a_v^{(k-1)}} da_v^{(k-1)} + \sum_{u=0}^{N^{(k)}-1} \Delta_u^{(k)} \frac{\partial z_u^{(k)}}{\partial b_u^{(k)}} db_u^{(L)}$$

$$\Rightarrow dC = \sum_{v=0}^{N^{(k-1)}-1} \sum_{u=0}^{N^{(k)}-1} \Delta_u^{(k)} \frac{\partial z_u^{(k)}}{\partial \tilde{w}_{vu}^{(k)}} d\tilde{w}_{vu}^{(k)} + \sum_{u=0}^{N^{(k)}-1} \Delta_u^{(k)} \frac{\partial z_u^{(k)}}{\partial b_u^{(k)}} db_u^{(L)} \\ + \sum_{v=0}^{N^{(k-1)}-1} \sum_{u=0}^{N^{(k)}-1} \Delta_u^{(k)} \frac{\partial z_u^{(k)}}{\partial a_v^{(k-1)}} da_v^{(k-1)}$$

(réorg. et inv. ordre sommation)

$$\Rightarrow dC = \sum_{v=0}^{N^{(k-1)}-1} \sum_{u=0}^{N^{(k)}-1} \Delta_u^{(k)} \frac{\partial z_u^{(k)}}{\partial \tilde{w}_{vu}^{(k)}} d\tilde{w}_{vu}^{(k)} + \sum_{u=0}^{N^{(k)}-1} \Delta_u^{(k)} \frac{\partial z_u^{(k)}}{\partial b_u^{(k)}} db_u^{(L)} \\ + \sum_{v=0}^{N^{(k-1)}-1} \Delta_v^{z(k)} da_v^{(k-1)}$$

$$(\Delta_v^{z(k)} = \sum_{u=0}^{N^{(k)}-1} \Delta_u^{(k)} \frac{\partial z_u^{(k)}}{\partial a_v^{(k-1)}})$$

$$\Rightarrow dC = \sum_{v=0}^{N^{(k-1)}-1} \sum_{u=0}^{N^{(k)}-1} \frac{\partial C}{\partial \tilde{w}_{vu}^{(k)}} d\tilde{w}_{vu}^{(k)} + \sum_{u=0}^{N^{(k)}-1} \frac{\partial C}{\partial b_u^{(k)}} db_u^{(L)} \\ + \sum_{v=0}^{N^{(k-1)}-1} \Delta_v^{z(k)} da_v^{(k-1)} \quad \blacksquare$$

$$\left(\frac{\partial C}{\partial \tilde{w}_{vu}^{(k)}} = \Delta_u^{(k)} \frac{\partial z_u^{(k)}}{\partial \tilde{w}_{vu}^{(k)}} \text{ et } \frac{\partial C}{\partial b_u^{(k)}} = \Delta_u^{(k)} \frac{\partial z_u^{(k)}}{\partial b_u^{(k)}} \right)$$

Le tableau récapitulatif des équations de la propagation du gradient

Afin d'évaluer les dérivées partielles

$$\frac{\partial C}{\partial \tilde{w}_{vu}^{(k)}} \quad \text{et} \quad \frac{\partial C}{\partial b_u^{(k)}}$$

de la fonction d'erreur par rapport à tous les poids et biais de toutes les couches dans une séquence d'activation tel que

$$a^{(*)} \rightarrow z^{(0)} \rightarrow a^{(0)} \rightarrow z^{(1)} \rightarrow a^{(1)} \rightarrow \dots \rightarrow z^{(k)} \rightarrow a^{(k)} \rightarrow z^{(k+1)} \rightarrow a^{(k+1)} \rightarrow \dots \rightarrow z^{(L)} \rightarrow a^{(L)},$$

nous devons calculer efficacement la différentielle de la fonction d'erreur.

En exploitant le concept de propagation de l'erreur $\Delta_u^{f(k)}$ de chacune des fonctions du réseau dans l'ordre inverse de l'activation, nous obtenons la séquence

$$C \rightarrow a^{(L)} \rightarrow z^{(L)} \rightarrow \dots \rightarrow a^{(k+1)} \rightarrow z^{(k+1)} \rightarrow a^{(k)} \rightarrow z^{(k)} \rightarrow \dots \rightarrow a^{(1)} \rightarrow z^{(1)} \rightarrow a^{(0)} \rightarrow z^{(0)} \rightarrow \Delta^{(*)}$$

les équations suivantes :

Propagation de l'erreur $\Delta_u^{f(k)}$			
Propagation de niveau $n = 0$	Propagation : $C^{(L)}$ tel $C^{(L)} \rightarrow a^{(L)}$	Propagation : $a^{(L)}$ tel $a^{(L)} \rightarrow z^{(L)}$	Indice des neurones
$L \rightarrow L-1$ (propagation de couche)	$\Delta_u^{C^{(L)}} = \frac{\partial C}{\partial a_u^{(L)}}$	$\Delta_u^{a^{(L)}} = \Delta_u^{C^{(L)}} \frac{\partial a_u^{(L)}}{\partial z_u^{(L)}}$	$u \in [0, N^{(L)} - 1]$
	P.S La fonction d'erreur possède le même nombre de neurones que la couche L	$\Delta_u^{a-full(L)} = \sum_{p=0}^{N^{(L)}-1} \Delta_p^{C^{(L)}} \frac{\partial a_p^{(L)}}{\partial z_u^{(L)}}$	
Propagation de niveau n où $k = L-n$	Propagation : $z^{(k)}$ tel $z^{(k)} \rightarrow a^{(k-1)}$	Propagation : $a^{(k)}$ tel $a^{(k)} \rightarrow z^{(k)}$	Indice des neurones
$k \rightarrow k-1$ (propagation de couche)	$\Delta_v^{z^{(k)}} = \sum_{u=0}^{N^{(k)}-1} \Delta_u^{a^{(k)}} \frac{\partial z_u^{(k)}}{\partial a_v^{(k-1)}}$	$\Delta_u^{a^{(k)}} = \Delta_u^{z^{(k+1)}} \frac{\partial a_u^{(k)}}{\partial z_u^{(k)}}$	$u \in [0, N^{(k)} - 1]$
	P.S Puisqu'il y a changement de couche, il y a changement du nombre de neurones passant de u à v .	$\Delta_u^{a-full(k)} = \sum_{p=0}^{N^{(k)}-1} \Delta_p^{z^{(k+1)}} \frac{\partial a_p^{(k)}}{\partial z_u^{(k)}}$	$v \in [0, N^{(k-1)} - 1]$

Gradient de la fonction d'erreur par rapport à aux paramètres $\tilde{w}_{vu}^{(k)}$ et $b_u^{(k)}$ du réseau	
$\frac{\partial C}{\partial \tilde{w}_{vu}^{(k)}} = \Delta_u^{a^{(k)}} \frac{\partial z_u^{(k)}}{\partial \tilde{w}_{vu}^{(k)}}$	$\frac{\partial C}{\partial b_u^{(k)}} = \Delta_u^{a^{(k)}} \frac{\partial z_u^{(k)}}{\partial b_u^{(k)}}$

Appendice

Les notations

Voici une liste des notations utilisées lors de la présentation des concepts du réseau de neurone :

* : Indice de la couche d'entrée.

k : Indice d'une couche de neurones. La première couche début à $k = 0$.

L : Indice de la dernière couche du réseau de neurones ($k \in [0, L]$). Puisque la numérotation en informatique des couches débute à $k = 0$, un réseau aura alors $L + 1$ couches.

i, j : Indices utilisés pour désigner un neurone particulier d'une couche de neurones.

$a^{(k)}$: Les valeurs de tous les neurones de la couche k .

$a_i^{(k)}$: Valeur du neurone i de la couche k .

$N^{(k)}$: Le nombre de neurones de la couche k . Ainsi $a^{(k)} = \{a_0^{(k)}, a_1^{(k)}, \dots, a_i^{(k)}, \dots, a_{N^{(k)}-1}^{(k)}\}$

$z^{(k)}$: Les agrégation de de tous les neurones de la couche k .

$z_i^{(k)}$: L'agrégation du neurone i de la couche k . Ce calcul nécessite des valeurs de neurones $a_j^{(k-1)}$ d'une couche précédente $k - 1$, des poids $w_{ij}^{(k)}$ et d'un biais $b_i^{(k)}$.

$w^{(k)}$: Les poids de tous les neurones de la couche k .

$w_{ij}^{(k)}$: Le poids du neurone $a_j^{(k-1)}$ dans l'agrégation $z_i^{(k)}$ du neurone i de la couche k .

$b^{(k)}$: Les biais de tous les neurones de la couche k .

$b_i^{(k)}$: Le biais dans l'agrégation $z_i^{(k)}$ du neurone i de la couche k .

C : La fonction d'erreur.

n : Indice du niveau de profondeur dans la descente du gradient. On utilisera une notation tel que $a^{(L-n)}$ pour désigner la valeur des neurones en itérant à partir de la fin du réseau.

α : Le taux d'apprentissage $\alpha > 0$.