

## Projet - Réseau de neurones : DIY

L'objectif de ce projet est d'implémenter un réseau de neurones. L'implémentation est inspirée des anciennes versions de `pytorch` (en Lua, avant l'autograd que vous verrez l'année prochaine) et des implémentations analogues qui permettent d'avoir des réseaux génériques très modulaires. Chaque couche du réseau est vu comme un module et un réseau est constitué ainsi d'un ensemble de modules. En particulier, les fonctions d'activation sont aussi considérées comme des modules (cf Figure 1).

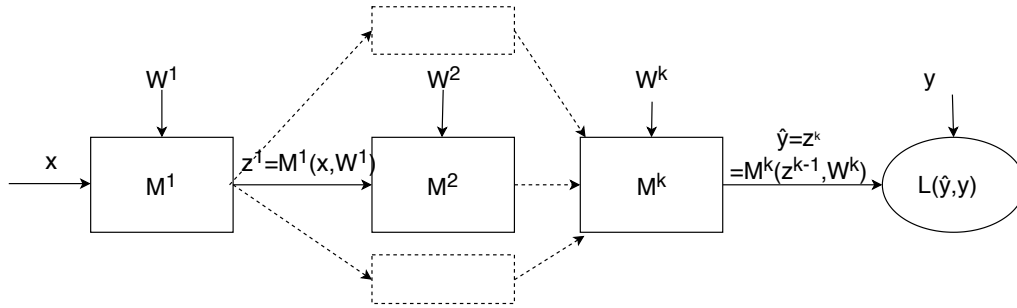


FIGURE 1 – Architecture module d'un réseau

Notons  $M^h(\mathbf{z}, \mathbf{W})$  le module de la couche  $h$  de paramètre  $\mathbf{W}$ ,  $\mathbf{z}^h = M^h(\mathbf{z}^{h-1}, \mathbf{W}^h)$  l'entrée de la couche  $h + 1$  (ou la sortie de la couche  $h$ ) et  $L(\mathbf{y}, \hat{\mathbf{y}})$  la fonction de coût. Pour pouvoir calculer la mise-à-jour des paramètres de chaque module  $h$ , on a besoin de calculer  $\nabla_{\mathbf{W}^h} L$ . Ce gradient est calculé par rétro-propagation et en utilisant la dérivation en chaîne. Il dépend de deux gradients :

- celui du module par rapport aux paramètres  $\nabla_{\mathbf{W}^h} M^h$  ; ce gradient est calculable sans connaître le reste du réseau, uniquement selon les caractéristiques du module ;
- celui de l'erreur par rapport aux sorties du module  $\nabla_{\mathbf{z}^h} L$  ; ce gradient est assimilable à l'erreur à corriger en rétro-propagation à la sortie du module, et est fourni par l'aval du réseau par induction (les modules  $M^{h+1}, M^{h+2}, \dots$ ). On note généralement les éléments de ce gradient  $\delta_j^h = \frac{\partial L}{\partial z_j^h}$

Ainsi pour un module  $M^h$  dont on "aplatit" les poids  $\mathbf{W}^h$  en une dimension  $(w_1^h, w_2^h, \dots, w_d^h)$ , on obtient les équations :

$$\frac{\partial L}{\partial w_i^h} = \sum_k \frac{\partial L}{\partial z_k^h} \frac{\partial z_k^h}{\partial w_i^h} = \sum_k \delta_k^h \frac{\partial z_k^h}{\partial w_i^h}, \text{ soit } \nabla_{\mathbf{w}^h} L = \begin{pmatrix} \frac{\partial z_1^h}{\partial w_1^h} & \frac{\partial z_2^h}{\partial w_1^h} & \dots \\ \frac{\partial z_1^h}{\partial w_2^h} & \ddots & \\ \vdots & & \end{pmatrix} \nabla_{\mathbf{z}^h} L \quad (1)$$

$$\delta_j^{h-1} = \frac{\partial L}{\partial z_j^{h-1}} = \sum_k \frac{\partial L}{\partial z_k^h} \frac{\partial z_k^h}{\partial z_j^{h-1}}, \text{ soit } \nabla_{\mathbf{z}^{h-1}} L = \begin{pmatrix} \frac{\partial z_1^h}{z_1^{h-1}} & \frac{\partial z_2^h}{z_1^{h-1}} & \dots \\ \frac{\partial z_2^h}{z_2^{h-1}} & \ddots & \dots \\ \vdots & & \end{pmatrix} \nabla_{\mathbf{z}^h} L \quad (2)$$

avec  $\frac{\partial z_k^h}{\partial z_i^{h-1}} = \frac{\partial M^h(\mathbf{z}^{h-1}, \mathbf{W}^h)_k}{\partial z_i^{h-1}}$ , la dérivée partielle de la  $k$ -ème sortie du module par rapport à la  $i$ -ème entrée.

Ainsi, pour pouvoir utiliser la back-propagation, il suffit que chaque module puisse calculer sa dérivée par rapport à ses paramètres (utilisée dans l'équation 1) et sa dérivée par rapport à ses entrées (utilisée dans l'équation 2).

## Classe (abstraite) Module

Cette section introduit le fonctionnement général de la librairie que vous allez développer. Elle est centrée autour de la classe abstraite `Module` qui représente un module générique du réseau de neurones. Le squelette vous est fourni dans le code source.

La classe `Module` contient :

- une variable `_parameters` qui stocke les paramètres du module lorsqu'il en a (la matrice de poids par exemple pour un module linéaire) ;
- une méthode `forward(data)` qui permet de calculer les sorties du module pour les entrées passées en paramètre ;
- une variable `gradient` qui permet d'accumuler le gradient calculé ;
- une méthode `zero_grad()` qui permet de réinitialiser à 0 le gradient ;
- une méthode `backward_update_gradient(input,delta)` qui permet de calculer le gradient du coût par rapport aux paramètres et l'additionner à la variable `_gradient` - en fonction de l'entrée `input` et des  $\delta$  de la couche suivante `delta` ;
- une méthode `backward_delta(input,delta)` qui permet de calculer le gradient du coût par rapport aux entrées en fonction de l'entrée `input` et des deltas de la couche suivante `delta` ;
- une méthode `update_parameters(gradient_step)` qui met à jour les paramètres du module selon le gradient accumulé jusqu'à son appel avec un pas de `gradient_step`.

Lorsque plusieurs modules sont mis en série, il suffit ainsi pour la passe forward d'appeler successivement les fonctions `forward` de chaque module avec comme entrée la sortie du précédent. Pour la passe backward, le dernier module calcule le gradient par rapport à ses paramètres et les deltas qu'il doit rétro-propager (à partir des deltas du loss) ; puis en parcourant en sens inverse le réseau, chaque module répète la même opération : le calcul de la mise à jour de son gradient (`backward_update_gradient`) et le delta qu'il doit transmettre à la couche précédente (`backward_delta`).

Remarquez que les paramètres ne sont pas mis tout de suite à jour en fonction du gradient : celui-ci est d'abord accumulé dans la variable `_gradient` et c'est uniquement lors de l'appel explicite à `backward_update_gradient` qui provoque la mise-à-jour des paramètres. Cela rend plus flexible l'utilisation des modules (plusieurs passes de backward peuvent être calculées avant de mettre à jour les paramètres du fait de l'additivité du gradient).

La classe `Loss` est plus simple : elle ne contient que deux méthodes :

- une fonction `forward(y,yhat)` qui permet de calculer le coût en fonction des deux entrées
- une fonction `backward(y,yhat)` qui permet de calculer le gradient du coût par rapport `yhat`.

Tout au long de votre implémentation, vous veillerez à réfléchir précisément à la taille des entrées et des sorties de chaque méthode. Il est conseillé d'utiliser l'instruction `assert condition` pour vous assurer que les paramètres que vous passez en entrée sont de la bonne taille. Par ailleurs, votre implémentation devra pouvoir traiter à chaque fois un `batch` d'exemples et non pas un seul exemple à la fois : ainsi la méthode `forward` d'une module linéaire devra pouvoir prendre en entrée une matrice de taille  $batch \times d$  ( $d$  la dimension des entrées).

## Mon premier est ... linéaire !

Pour cette première étape, vous allez coder les deux classes dont vous avez besoin pour réaliser une régression linéaire :

- une fonction de coût `MSELoss` dont la méthode `forward(y,yhat)` doit rendre  $\|y - \hat{y}\|^2$  ; Attention à la généralité de votre implémentation : la supervision `y` et la prédiction `yhat` sont des matrices de taille  $batch \times d$  (chaque supervision peut être un vecteur de taille  $d$ , pas seulement un scalaire comme dans le cas de la régression univariée). La fonction doit rendre un vecteur de dimension `batch` (le nombre d'exemples).

- un module `Linear(input,output)` qui représente une couche linéaire avec `input` entrées et `output` sorties. La méthode `forward` prend donc une matrice de taille  $batch \times input$  et produit une sortie  $batch \times output$ .

N'oubliez pas de coder toutes les fonctions de ces deux modules! Une fois l'implémentation réalisée, testez-la sur des données quelconques en réalisant une boucle d'apprentissage par descente de gradient pour optimiser votre premier réseau.

### Mon second est ... non-linéaire !

Implémentez le module `TanH` qui permet d'appliquer une tangente hyperbolique aux entrées et le module `Sigmoïde` qui permet d'appliquer une sigmoïde aux entrées. N'oubliez pas que les modules de transformation héritent de la classe `Module` et donc doivent implémenter les fonctions `backward_update_gradient`, `backward_delta` et `update_parameters` même si le module n'a pas de paramètre!

Testez votre implémentation en réalisant un réseau à deux couches linéaires avec une activation tangente entre les deux couches et une activation sigmoïde à la sortie. Vous utiliserez des données d'un problème de classification binaire en considérant 1 et 0 comme classes positive et négative.

### Mon troisième est un encapsulage

En réalisant le réseau à deux couches précédents, vous remarquez que les opérations de chaînage entre modules sont répétitives lors de la descente de gradient - que ce soit pour la passe forward ou backward - et qu'il sera fastidieux de les écrire pour un grand nombre de modules. Implémenter une classe `Sequential` qui permet d'ajouter des modules en série et qui automatise les procédures de forward et backward quel que soit le nombre de modules mis à la suite.

Après l'avoir testé, vous pouvez implémenter une classe `Optim(net,loss,eps)` pour condenser une itération de gradient : elle prend dans son constructeur un réseau `net`, une fonction de coût `loss` et un pas `eps`. Elle contient une seule méthode `step(batch_x,batch_y)` qui calcule la sortie du réseau sur `batch_x`, calcule le coût par rapport aux labels `batch_y`, exécute la passe backward et met à jour les paramètres du réseau.

Vous pouvez également implémenter une fonction `SGD` qui prend en entrée entre autre un réseau, un jeu de données, une taille de batch et un nombre d'itération et s'occupe du découpage en batch du jeu de données et de l'apprentissage du réseau pendant le nombre d'itérations spécifié.

### Mon quatrième est multi-classe

Le multi-classe utilise en sortie du réseau une dimension par classe pour dénoter la probabilité de chaque classe. Le vecteur de supervision est un encodage one-hot : un vecteur rempli de 0 sauf à l'index de la bonne classe qui prend la valeur 1. On utilise un Softmax que l'on introduit à la dernière couche pour transformer les entrées en distribution de probabilités grâce à la normalisation effectuée. On peut utiliser la MSE comme coût, cependant elle donne des résultats décevants car elle a tendance à trop "moyenner" les erreurs et ne pas pousser les sorties vers 0 ou 1. Il est mieux d'utiliser des coûts adaptés aux distributions de probabilités comme la cross entropie.

Pour pouvoir faire du multi-classe, nous avons donc besoin :

- d'une transformation `Softmax` qui permet d'appliquer un soft-max aux entrées :  $\text{softmax}(\mathbf{z}) = \left( \frac{e^{z_1}}{\sum_k e^{z_k}}, \frac{e^{z_2}}{\sum_k e^{z_k}}, \dots \right)$
- d'un coût cross-entropique : pour  $y$  l'indice de la classe à prédire et  $\hat{\mathbf{y}}$  le vecteur de prédiction, le coût cross-entropique (équivalent à un maximum de vraisemblance) est  $\text{CE}(y, \hat{\mathbf{y}}) = -\hat{\mathbf{y}}_y$  ;
- il est habituel de combiner les deux ensembles afin d'éviter des instabilités numériques. Dans ce cas, on enchaîne un `Softmax` passé au logarithme (`logSoftMax`) et un coût cross entropique. Le

coût s'écrit alors :  $CE(y, \hat{y}) = -\log \frac{e^{\hat{y}_y}}{\sum_{i=1}^K e^{\hat{y}_i}} = -\hat{y}_y + \log \sum_{i=1}^K e^{\hat{y}_i}$ .

Testez sur le jeu de données des chiffres manuscrits par exemple.

## Mon cinquième se comprime ( et s'expérimente beaucoup)

Un auto-encodeur est un réseau de neurones dont l'objectif est d'apprendre un encodage des données dans le but généralement de réduire les dimensions. Il s'agit d'apprentissage non-supervisé : il n'y a pas de classes associées à chaque exemple lors de l'entraînement. La réduction de dimension dans ce cas à de multiples applications : débruitage d'image, reconstruction de parties cachées, clustering, visualisation ... On peut également se servir de l'exemple encodé comme nouvelle description de l'exemple pour ensuite faire de l'apprentissage supervisé classique.

Concrètement, un auto-encodeur est formé de deux parties :

- un *encodeur* qui prend en entrée une donnée et la transforme dans un espace de plus petite dimension (à l'aide par exemple de couches linéaires successives de plus en plus petites, alternées avec des fonctions d'activation non linéaires) pour obtenir le *code* (ou représentation latente) de l'exemple ;
- un *décodeur* qui prend en entrée le code d'un exemple et généralement avec un réseau symétrique à l'encodeur décode l'exemple vers sa représentation initiale.

Un exemple d'architecture d'auto-encodeur :

- Encodage :  $\text{Linear}(256, 100) \rightarrow \text{TanH}() \rightarrow \text{Linear}(100, 10) \rightarrow \text{TanH}()$
- Décodage :  $\text{Linear}(10, 100) \rightarrow \text{TanH}() \rightarrow \text{Linear}(100, 256) \rightarrow \text{Sigmoïde}()$

Dans cet exemple, une compression dans un espace de dimension 10 est réalisée à partir d'une entrée en taille 256. La sigmoïde finale permet de s'assurer que les sorties sont entre 0 et 1 (utile si c'est le cas pour la donnée d'entrée également).

L'entraînement se fait sur un coût de reconstruction :  $L(X, \hat{X})$ , avec  $\hat{X} = \text{decodeur}(\text{encodeur}(X))$  la donnée reconstruite. Le coût de reconstruction peut être par exemple un coût aux moindres carrés ou - plus performant - une cross entropie binaire :  $\text{BCE}(y, \hat{y}) = -(y * \log(\hat{y}) + (1 - y)\log(1 - \hat{y}))$ .

La cross-entropie binaire pour une sortie  $y$  entre 0 et 1 permet d'être plus "abrupte" que la MSE sur les valeurs de sorties, i.e. de pousser les valeurs vers 0 ou 1 plutôt que vers un moyennage de valeurs comme le fait la MSE (quel rapport voyez vous avec la vraisemblance ?). Elle est appelée binaire car la supervision  $y$  dans ce cas est censée être binaire (soit 0 soit 1), dans les faits elle peut être utilisée sur des valeurs réelles compris entre 0 et 1. En pratique, on utilise également un seuillage (ou clip) du  $\log$  pour ne pas avoir des valeurs infinies lorsque  $\hat{y}$  ou  $1 - \hat{y}$  tendent vers 0, en retournant par exemple  $\max(-100, \log(\hat{y}))$ .

Il est habituel d'utiliser pour le décodeur les matrices de poids transposées de l'encodeur afin d'assurer une régularisation (ce qu'on appelle un *partage de poids* dans le réseau), ce n'est pas demandé dans la suite.

Expérimentez ce type d'architecture avec les modules que vous avez développés en faisant varier les hyperparamètres (nombre de couches, taille de l'espace latent, fonctions d'activation, ...). Vous pouvez en particulier :

- visualiser les images reconstruites après une forte compression
- étudier le clustering induit dans l'espace latent (avec du  $k$ -means par exemple) et en utilisant les classes des exemples étudier la performance des représentations apprises
- visualiser les représentations obtenues dans un espace 2D ou 3D, soit directement si votre dimension latente est de 2 ou 3, soit en utilisant un algorithme tel que  $t$ -sne pour projeter vos représentations en 2D
- étudier les performances en débruitage : bruitez les données et observer l'erreur entre la donnée débruitée par le réseau et la donnée originale. Vous pouvez également utiliser des données bruitées lors de l'apprentissage pour améliorer les performances.
- étudier les performances en classification en utilisant comme représentation la représentation

latente construite : vous utilisez dans ce cas l'auto-encodeur comme un pré-traitement de vos données.

Il n'est pas demandé à ce que vous traitiez tous les sujets : choisissez une piste qui vous intéresse et montez un plan expérimental cohérent. Présentez l'objectif et les résultats auxquels vous parvenez.

## Mon sixième se convole

La couche convolutionnelle est le standard en classification d'image. Elle permet d'appliquer un même *filtre* sur différentes positions de l'image et de sortir une valeur par position. Les sorties associées à un filtre sont une nouvelle image (de taille différente en fonction du nombre de position où l'opération de convolution a été effectuée) qui est appelée *feature map*. Vous pouvez vous référer à ce résumé pour de belles animations (ou plein d'autres sites il y en a à foison).

Pour simplifier l'implémentation, nous nous intéresserons uniquement aux convolutions en 1 dimension. Nous supposons donc notre image d'entrée *aplatit* dans une matrice de taille  $d \times C$ ,  $d = \text{largeur} \times \text{hauteur}$  comme dans les parties précédentes. Le  $C$  correspond au nombre de canaux : pour une image en noir et blanc, il n'y en a qu'un, pour une image en couleur il y en a typiquement 3 (le codage le plus courant étant le codage RGB - rouge, vert, bleu). La notation  $x_{i,c}$  désignera la valeur du  $i$ -ème pixel du canal  $c$ .

Un filtre de convolution est un opérateur  $m_W$  caractérisé par sa taille (*kernel size*, notée  $ksize$  par la suite). Ses paramètres  $W$  sont de tailles  $kernel\ size * C$ . Pour une position  $i$  dans l'image, l'opérateur renvoie la sortie  $o_i = \sum_{j=1}^{ksize} \sum_{c=1}^C w_{j,c} x_{i+j-1,c}$ . Ainsi le filtre correspond à un opérateur linéaire sur une petite partie de l'image - une fenêtre - située entre le pixel  $i$  et le pixel  $i + ksize - 1$ .

Le filtre peut être appliqué successivement sur tous les pixels d'une image. On obtient alors une sortie de taille  $d - ksize + 1$  (en effet on ne peut pas appliquer le filtre sur les dernières positions de l'image, le filtre déborderait des bords de l'image). Afin de réduire le temps de traitement, de produire des sorties plus compactes et d'éviter une trop grande redondance, un filtre a également un paramètre appelé *stride* : il s'agit du déplacement élémentaire effectué lors de l'application du filtre. Au lieu d'appliquer le filtre à toutes les positions, le filtre sera appliqué tous les *stride* positions. La sortie sera alors  $o_1, o_{1+stride}, o_{1+2*stride} \dots$  et de taille  $\lfloor (d - ksize) / stride \rfloor + 1$  (la partie entière des positions possibles divisée par le stride incrémentée de 1). Un autre paramètre peut être pris en considération, le *padding* qui permet d'appliquer le filtre sur les bord de l'image, nous ne l'étudierons pas dans ce sujet. Pour résumer : un filtre est défini par sa taille **ksize** et son *stride*, ses paramètres de taille  $ksize \times C$  et lorsqu'il est appliqué à une image de taille  $d \times C$  il sort un vecteur de taille  $d_{out} = \lfloor (d - ksize) / stride \rfloor + 1$ .

La sortie  $o_i$  du vecteur est  $\sum_{j=1}^{ksize} \sum_{c=1}^C w_{i*stride+j-1,c}$ .

Une couche convolutionnelle est composée de plusieurs filtres de mêmes paramètres **ksize** et **stride** - ils produisent des vecteurs de même dimension en sortie - mais bien sûr de paramètres  $W$  différents. Le nombre de filtres est souvent noté  $C_{out}$  - nombre de canaux de sorties. Les vecteurs de sortie sont concaténées pour former la *feature map* de dimension  $d_{out} \times C_{out}$ , considérée comme nouvelle image et qui peut être à son tour passée dans une nouvelle couche convolutionnelle.

L'autre type de couche utilisé en alternance avec les couches convolutionnelles est la couche de *pooling* qui permet de réduire la taille des feature maps considérées en effectuant un sous-échantillonnage. Elle fonctionne comme un filtre mais au lieu d'appliquer une transformation linéaire, elle calcule une valeur agrégée des valeurs présentes dans la fenêtre considérée : la moyenne, le maximum pour les cas les plus courants. Elle a donc également un paramètre **ksize** qui définit la taille de fenêtre considérée et un paramètre **stride** qui définit son déplacement dans l'image. Elle ne dispose pas d'autres paramètres. Dans le cas d'un max-pooling par exemple, la sortie  $o_i$  renvoyée à une position  $i$  est :

$$\max_{j=1..ksize, c=1..C} x_{i+j-1,c}$$

Il reste deux éléments à introduire pour pouvoir faire un réseau convolutionnel. Le premier est la fonction **ReLU** - (Rectified Linear Unit) que l'on préfère comme fonction d'activation aux autres usuelles :  $ReLU(x) = \max(0, x)$ . Elle permet entre autre d'éviter les problèmes de gradient vanishing dans les réseaux profonds. Enfin, pour aboutir à une classification finale, on utilise une couche linéaire à partir

de la sortie de la dernière couche convolutionnelle (ou de pooling). Pour cela, il faut aplatir la sortie de cette couche. Il faut donc un module `Flatten` qui permet de transformer une entrée 2D  $d_{out} \times C$  en une sortie de taille  $d_{out} * C$ .

En résumé, pour pouvoir réaliser un réseau convolutionnel, il faut coder :

- un module `Conv1D(k_size,chan_in,chan_out,stride)` qui contient une matrice de paramètres de taille  $(k\_size,chan\_in,chan\_out)$  (soit  $chan\_out$  filtres de taille  $(k\_size,chan\_in)$ ); sa méthode `forward` prend en entrée un batch de taille  $(batch,length,chan\_in)$  et sort une matrice de taille  $(batch,(length-k\_size)/stride +1,chan\_out)$ .
- un module `MaxPool1D(k_size,stride)` sans paramètres à apprendre. La méthode `forward` prend en entrée un batch de taille  $(batch,length,chan\_in)$  et sort une matrice de taille  $(batch,(length-k\_size)/stride +1,chan\_in)$ .
- un module `Flatten` qui prend en entrée un batch de taille  $(batch,length,chan\_in)$  et renvoie une matrice  $(batch, length * chan\_in)$ ;
- une fonction d'activation `ReLU`.

Implémentez les différents modules et expérimentez différentes architectures pour la reconnaissance des chiffres. À titre indicatif, un réseau de type `Conv1D(3,1,32) → MaxPool1D(2,2) → Flatten()` `→ Linear(4064,100) → ReLU() → Linear(100,10)` fait de très bonnes performances.

## Mon tout s'améliore

Les convolutions 1D implémentées permettent de reconnaître des motifs horizontaux dans les images. À partir de vos modules, il est très facile de reconnaître également des motifs verticaux : il suffit pour cela de transposer l'image. Vous pouvez donc implémenter un réseau qui reconnaît des motifs verticaux et horizontaux en concaténant la sortie de deux couches convolutionnelles. Vous pouvez également considérer du Average Pooling (qui fait la moyenne sur la fenêtre temporelle) plutôt que du Max Pooling. Les plus courageux peuvent enfin implémenter les "vraies" convolutions 2D. Ce n'est que quelques pistes possibles d'améliorations, cette partie est laissée libre à vos envies selon votre temps disponible (d'autres données images, d'autres types de données par exemple temporelles, ...).

## Quelques astuces implémentatoires

- pour transformer un vecteur d'indices `y` de classes en one-hot : `onehot = np.zeros((y.size,10)); onehot[np.arange(y.size),y]=1`
- pour la partie convolutionnelle, vous aurez à utiliser `np.newaxis` qui permet d'ajouter une dimension à un tableau (équivalent à un `.reshape(1, ...)`).
- une initialisation aléatoire des poids entre -1 et 1 des couches convolutionnelles mènera potentiellement à des nombres trop grands en sortie (et donc à un NaN au passage de l'exponentiel). Si cela se produit, modifiez l'initialisation par un facteur  $10^{-1}$  ou plus.
- pour le MaxPool, vous aurez besoin pour le gradient d'affecter des valeurs que aux endroits où le max est atteint (et des zéros partout ailleurs). Une façon de faire est d'utiliser la même astuce que pour le one hot : si `idx` contient les indices de l'argmax sous forme aplati, soit `res` la matrice résultat de taille `nb_batch,length,chan` : `res[np.repeat(range(nb_batch),chan),idx,list(range(chan))*nb_batch]` permet de toucher les indices voulues.
- de manière générale, n'oubliez pas que le gradient est de la même taille que les paramètres et que `backward_delta` prend un delta de même taille que la sortie du `forward` du module.
- essayez de faire le moins de boucles possibles ...