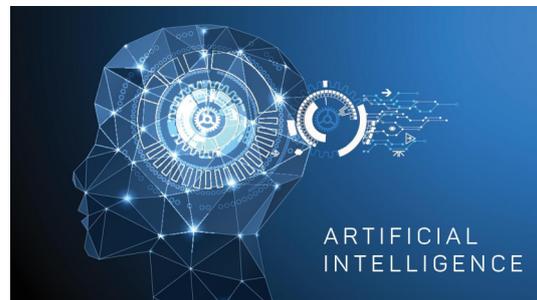


# Chapitre 8 : Réseaux de neurones en pratique pour l'apprentissage automatique



Pr. Mustapha BOURAHLA, Département  
d'Informatique, Université de M'Sila, Contact :  
[mustapha.bourahla@univ-msila.dz](mailto:mustapha.bourahla@univ-msila.dz)

# Table des matières



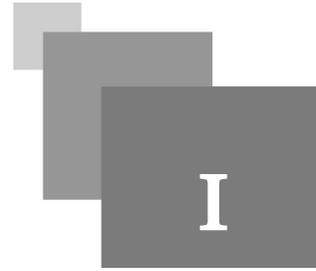
<b>Introduction</b>	3
<b>I - Introduction</b>	4
<b>II - Couche par Couche</b>	5
<b>III - Descente de Gradient</b>	6
<b>IV - Classe abstraite : Layer</b>	8
<b>V - Fully Connected Layer</b>	9
<b>VI - Couche d'activation</b>	13
<b>VII - Couche convolution</b>	15
<b>VIII - Flatten layer</b>	19
<b>IX - Fonction de perte</b>	21
<b>X - Classe Network</b>	22
<b>XI - Construire un réseau de neurone</b>	24
<b>XII - exemple MNIST</b>	26
<b>XIII - Le deep learning</b>	28
<b>XIV - Exercice :</b>	30
<b>XV - Exercice :</b>	33
<b>Conclusion</b>	35

# Introduction



Ce chapitre vous montre comment programmer un réseau de neurones en Python pour réaliser un apprentissage automatique.

# Introduction



On supposera à partir d'ici que vous avez déjà quelques connaissances fondamentales sur le modèle du neurone artificiel et les réseaux de neurones. Dans ce chapitre nous présentons l'algorithme des gradients descendants. Encore une fois, le but ici n'est pas d'expliquer les différentes applications possibles du machine learning, mais plutôt comment implémenter ces algorithmes.

# Couche par Couche

II

Gardons à l'esprit la démarche globale du machine learning :

1. Donner une entrée au modèle.
2. Propager cette entrée à travers le réseau de neurones jusqu'à récupérer la sortie.
3. Une fois la sortie obtenue, nous pouvons la comparer à la sortie voulue et donc calculer une erreur.
4. On ajuste les paramètres du modèle pour diminuer l'erreur précédemment calculée. Pour cela on soustrait à chaque paramètre la dérivée de l'erreur par rapport à lui-même (gradient descendant).
5. On recommence à l'étape 1.

L'étape la plus importante est la 4ième. Nous voulons être capable de créer autant de couches que l'on veut, de n'importe quel type, et d'utiliser n'importe quelle fonction d'activation. Seulement, en changeant l'architecture du réseau de neurones, on change également la formule littérale du calcul de la dérivée de l'erreur par rapport aux paramètres.

Le but est donc de faire une implémentation qui fait abstraction de l'architecture du modèle (comme dans Keras). Pour cela, nous devons implémenter chaque couche séparément.

## *Ce que chaque couche doit faire*

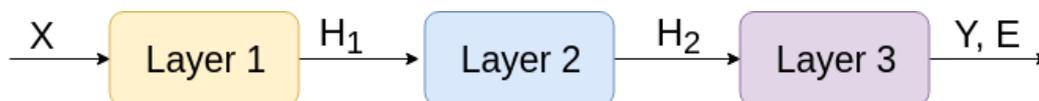
Quelle que soit la couche que nous codons (fully connected, convolutional, maxpooling, dropout, etc.), il y aura toujours au moins deux éléments fondamentaux : une entrée et une sortie.



## *Passé avant — forward propagation*

Nous pouvons dès lors préciser une propriété importante : la sortie d'une couche est l'entrée de la couche suivante.

Cette partie est ce qu'on appelle la passe avant (forward propagation) : on propage l'entrée X (image, son, texte, etc.) dans le réseau de neurones jusqu'à obtenir la sortie Y. Puis, on observe une erreur E qu'il faut maintenant diminuer.



# Descente de Gradient

III

Fondamentalement, nous voulons changer un paramètre dans le réseau (appelez-le  $w$ ) afin que l'erreur totale  $E$  diminue. Il existe un moyen intelligent de le faire (sans changer le paramètre au hasard) qui est le suivant:

$$w \leftarrow w - \alpha \frac{\partial E}{\partial w}$$

Où  $\alpha$  est un paramètre dans l'intervalle  $[0,1]$  que nous fixons et qui est appelé le taux d'apprentissage. Quoi qu'il en soit, l'important ici est  $\partial E/\partial w$  (la dérivée de  $E$  par rapport à  $w$ ). Nous devons être en mesure de trouver la valeur de cette expression pour n'importe quel paramètre du réseau, quelle que soit son architecture.

*Passer arrière — backward propagation*

Supposons que l'on donne à une couche la dérivée de l'erreur par rapport à sa sortie ( $\partial E/\partial Y$ ), alors elle doit être capable de donner la dérivée de l'erreur par rapport à son entrée ( $\partial E/\partial X$ ).

$$\frac{\partial E}{\partial X} \leftarrow \boxed{\text{layer}} \leftarrow \frac{\partial E}{\partial Y}$$

L'erreur  $E$  est un scalaire (un nombre), et  $X$  et  $Y$  sont des matrices. La notation ci-dessus (abusive) signifie ceci:

$$\frac{\partial E}{\partial X} = \left[ \frac{\partial E}{\partial x_1} \quad \frac{\partial E}{\partial x_2} \quad \cdots \quad \frac{\partial E}{\partial x_i} \right]$$

$$\frac{\partial E}{\partial Y} = \left[ \frac{\partial E}{\partial y_1} \quad \frac{\partial E}{\partial y_2} \quad \cdots \quad \frac{\partial E}{\partial y_j} \right]$$

Laissons de côté  $\partial E/\partial X$  pour l'instant et concentrons-nous sur  $\partial E/\partial Y$ . Si une couche a accès à  $\partial E/\partial Y$  où  $Y$  est sa propre sortie, alors nous pouvons très facilement calculer la dérivée de l'erreur par rapport à ses paramètres  $\partial E/\partial W$  (pour l'étape de l'ajustement), et cela, indépendamment de l'architecture globale du réseau de neurones! Il suffit d'utiliser la règle de dérivation des fonctions composées :

$$\frac{\partial E}{\partial w} = \sum_j \frac{\partial E}{\partial y_j} \frac{\partial y_j}{\partial w}$$

Etant donnée  $\partial E/\partial Y$  nous pouvons donc calculer  $\partial E/\partial W$ , et donc ajuster les paramètres de la couche!

### Pourquoi avons-nous besoin de $\partial E/\partial X$ ?

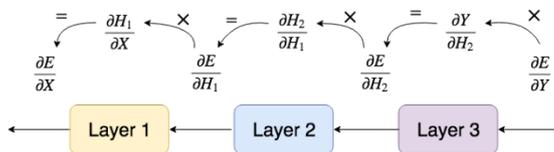
N'oubliez pas, la sortie d'une couche est l'entrée de la couche suivante. Donc  $\partial E/\partial X$  pour une couche sera  $\partial E/\partial Y$  pour la couche précédente! Une fois munit de son propre  $\partial E/\partial Y$ , la couche précédente pourra à son tour ajuster ses paramètres. Pour calculer  $\partial E/\partial X$  on utilise encore une fois la règle de dérivation des fonctions composées :

$$\frac{\partial E}{\partial x_i} = \sum_j \frac{\partial E}{\partial y_j} \frac{\partial y_j}{\partial x_i}$$

Cette astuce est la clé de compréhension de la backward propagation! Après cela, nous pourrons programmer un réseau de neurones convolutif en un rien de temps.

### Un super diagramme

C'est ce que j'ai décrit plus tôt. La couche 3 va mettre à jour ses paramètres en utilisant  $\partial E/\partial Y$ , puis va passer  $\partial E/\partial H_2$  à la couche précédente, qui est son propre " $\partial E/\partial Y$ ". La couche 2 va alors faire de même, et ainsi de suite.



Ça peut paraître abstrait maintenant mais deviendra très clair part la suite. Nous pouvons dès lors créer notre première classe en Python qui sera une classe abstraite représentant une couche.

# Classe abstraite : Layer

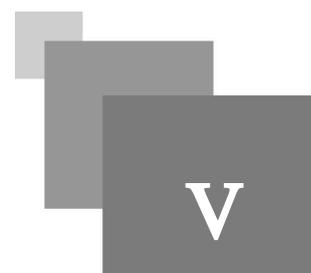
## IV

### Classe abstraite : Layer

- La classe abstraite Layer, dont les autres couches hériteront, contient les caractéristiques communes à toutes les couches :
  1. une entrée,
  2. une sortie,
  3. une fonction qui fait la passe avant, et une pour la passe arrière.
- Comme vous pouvez le constater, il existe un paramètre supplémentaire dans `backward_propagation` qu'il n'est pas mentionné, c'est le `learning_rate`.
- Ce paramètre devrait être une politique de mise à jour, ou un optimiseur, comme ils l'appellent dans Keras, mais par souci de simplicité, nous allons simplement passer le learning rate et mettre à jour nos paramètres en utilisant la descente de gradient.

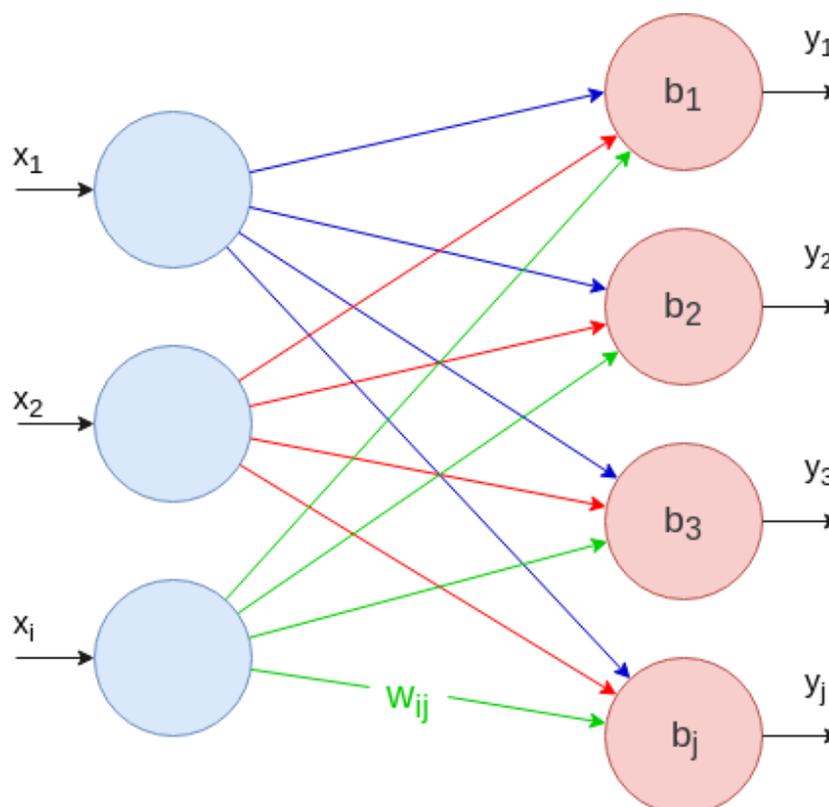
```
1 # Base class
2 class Layer:
3     def __init__(self):
4         self.input = None
5         self.output = None
6
7     # computes the output Y of a layer for a given input X
8     def forward_propagation(self, input):
9         raise NotImplementedError
10
11    # computes dE/dX for a given dE/dY (and update parameters if any)
12    def backward_propagation(self, output_error, learning_rate):
13        raise NotImplementedError
```

# Fully Connected Layer



## *Fully Connected Layer*

Définissons et implémentons maintenant le premier type de couche: Fully Connected Layer ou FC Layer. Les couches FC sont les couches les plus élémentaires car tous les neurones d'entrée sont connectés à tous les neurones de sortie.



## *Forward Propagation*

La valeur de chaque neurone de sortie peut être calculée comme suit :

$$y_j = b_j + \sum_i x_i w_{ij}$$

La notation matricielle permet de simplifier ce calcul :

$$X = [x_1 \quad \dots \quad x_i] \quad W = \begin{bmatrix} w_{11} & \dots & w_{1j} \\ \vdots & \ddots & \vdots \\ w_{i1} & \dots & w_{ij} \end{bmatrix} \quad B = [b_1 \quad \dots \quad b_j]$$

$$Y = XW + B$$

Nous en avons fini avec la passe avant. Traitons maintenant la passe arrière de la couche FC.

Notez que je n'utilise aucune fonction d'activation, c'est parce que nous allons l'implémenter dans une couche à part !

### Backward Propagation

Comme nous l'avons dit, supposons que nous ayons une matrice contenant la dérivée de l'erreur par rapport à la sortie de cette couche ( $\partial E / \partial Y$ ). Nous avons besoin de :

La dérivée de l'erreur par rapport aux paramètres ( $\partial E / \partial W$ ,  $\partial E / \partial B$ )

La dérivée de l'erreur par rapport à l'entrée ( $\partial E / \partial X$ )

Commençons par  $\partial E / \partial W$ . Cette matrice doit avoir la même taille que  $W$  :  $i \times j$  où  $i$  est le nombre de neurones d'entrée et  $j$  le nombre de neurones de sortie. Nous avons besoin d'une dérivée pour chaque paramètre :

$$\frac{\partial E}{\partial W} = \begin{bmatrix} \frac{\partial E}{\partial w_{11}} & \dots & \frac{\partial E}{\partial w_{1j}} \\ \vdots & \ddots & \vdots \\ \frac{\partial E}{\partial w_{i1}} & \dots & \frac{\partial E}{\partial w_{ij}} \end{bmatrix}$$

En utilisant la règle de dérivation des fonctions composées, on écrit :

$$\begin{aligned} \frac{\partial E}{\partial W} &= \begin{bmatrix} \frac{\partial E}{\partial y_1} & \dots & \frac{\partial E}{\partial y_j} \\ \vdots & \ddots & \vdots \\ \frac{\partial E}{\partial y_1} & \dots & \frac{\partial E}{\partial y_j} \end{bmatrix} \\ &= \begin{bmatrix} x_1 \\ \vdots \\ x_i \\ x_i \end{bmatrix} \begin{bmatrix} \frac{\partial E}{\partial y_1} & \dots & \frac{\partial E}{\partial y_j} \end{bmatrix} \\ &= X^T \frac{\partial E}{\partial Y} \end{aligned}$$

Nous avons notre première formule qui permet d'ajuster les poids ! Calculons à présent  $\partial E / \partial B$ .

$$\frac{\partial E}{\partial B} = \left[ \frac{\partial E}{\partial b_1} \quad \frac{\partial E}{\partial b_2} \quad \dots \quad \frac{\partial E}{\partial b_j} \right]$$

Encore une fois,  $\partial E / \partial B$  doit-être de la même dimension que  $B$ : une dérivée par biais.

$$\begin{aligned} \frac{\partial E}{\partial b_j} &= \frac{\partial E}{\partial y_1} \frac{\partial y_1}{\partial b_j} + \dots + \frac{\partial E}{\partial y_j} \frac{\partial y_j}{\partial b_j} \\ &= \frac{\partial E}{\partial y_j} \end{aligned}$$

D'où,

$$\begin{aligned} \frac{\partial E}{\partial B} &= \left[ \frac{\partial E}{\partial y_1} \quad \frac{\partial E}{\partial y_2} \quad \dots \quad \frac{\partial E}{\partial y_j} \right] \\ &= \frac{\partial E}{\partial Y} \end{aligned}$$

Maintenant que nous avons  $\partial E/\partial W$  et  $\partial E/\partial B$  nous pouvons ajuster tout les paramètres de cette couche de sorte à diminuer l'erreur! Il nous reste simplement à calculer  $\partial E/\partial X$  pour que la couche précédente puisse faire les mêmes calculs.

$$\frac{\partial E}{\partial X} = \left[ \frac{\partial E}{\partial x_1} \quad \frac{\partial E}{\partial x_2} \quad \dots \quad \frac{\partial E}{\partial x_i} \right]$$

Encore une fois, dérivation de fonctions composées :

$$\begin{aligned} \frac{\partial E}{\partial x_i} &= \frac{\partial E}{\partial y_1} \frac{\partial y_1}{\partial x_i} + \dots + \frac{\partial E}{\partial y_j} \frac{\partial y_j}{\partial x_i} \\ &= \frac{\partial E}{\partial y_1} w_{i1} + \dots + \frac{\partial E}{\partial y_j} w_{ij} \end{aligned}$$

Nous pouvons écrire la matrice entière :

$$\begin{aligned} \frac{\partial E}{\partial X} &= \left[ \left( \frac{\partial E}{\partial y_1} w_{11} + \dots + \frac{\partial E}{\partial y_j} w_{1j} \right) \quad \dots \quad \left( \frac{\partial E}{\partial y_1} w_{i1} + \dots + \frac{\partial E}{\partial y_j} w_{ij} \right) \right] \\ &= \begin{bmatrix} \frac{\partial E}{\partial y_1} & \dots & \frac{\partial E}{\partial y_j} \end{bmatrix} \begin{bmatrix} w_{11} & \dots & w_{1j} \\ \vdots & \ddots & \vdots \\ w_{i1} & \dots & w_{ij} \end{bmatrix} \\ &= \frac{\partial E}{\partial Y} W^t \end{aligned}$$

Et voilà! Nous avons obtenu ces trois formules fondamentale pour la couche FC!

$$\begin{aligned} \frac{\partial E}{\partial X} &= \frac{\partial E}{\partial Y} W^t \\ \frac{\partial E}{\partial W} &= X^t \frac{\partial E}{\partial Y} \\ \frac{\partial E}{\partial B} &= \frac{\partial E}{\partial Y} \end{aligned}$$

### Coder la couche FC

Créez un nouveau fichier qui contiendra le code suivant :

```

1 from layer import Layer
2 import numpy as np
3
4 # inherit from base class Layer
5 class FCLayer(Layer):
6     # input_size = number of input neurons
7     # output_size = number of output neurons
8     def __init__(self, input_size, output_size):
9         self.weights = np.random.rand(input_size, output_size) - 0.5
10        self.bias = np.random.rand(1, output_size) - 0.5
11
12    # returns output for a given input
13    def forward_propagation(self, input_data):
14        self.input = input_data
15        self.output = np.dot(self.input, self.weights) + self.bias
16        return self.output
17
18    # computes dE/dW, dE/dB for a given output_error=dE/dY. Returns input_error=dE/dX.
19    def backward_propagation(self, output_error, learning_rate):
20        input_error = np.dot(output_error, self.weights.T)
21        weights_error = np.dot(self.input.T, output_error)

```

```
22     # dBias = output_error
23
24     # update parameters
25     self.weights -= learning_rate * weights_error
26     self.bias -= learning_rate * output_error
27     return input_error
```

# Couche d'activation

VI

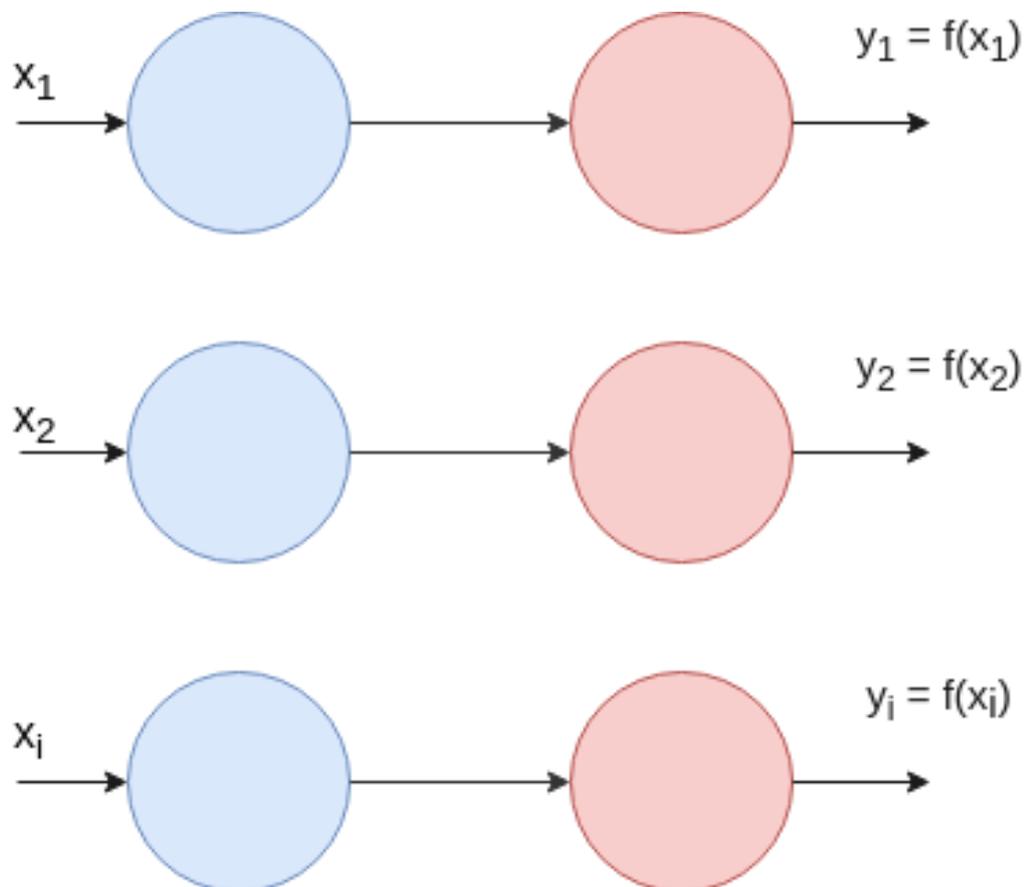
## *Couche d'activation*

Tous les calculs que nous avons faits jusqu'à présent étaient complètement linéaires. La machine n'apprendra rien avec ce genre de modèle. Nous devons ajouter une non-linéarité au modèle en appliquant des fonctions non linéaires à la sortie de certaines couches.

Nous devons maintenant refaire tout le processus pour ce nouveau type de couche !

Pas de soucis, ça va être bien plus rapide car il n'y a pas de paramètres entraînable. Il suffit de calculer  $\partial E / \partial X$ .

Nous appellerons  $f$  et  $f'$  la fonction d'activation et sa dérivée, respectivement.



## *Forward Propagation*

Comme vous le verrez, c'est assez simple. Pour une entrée  $X$  donnée, la sortie est simplement la fonction d'activation appliquée à chaque élément de  $X$ . Ce qui signifie que l'entrée et la sortie ont la même dimension.

$$Y = \begin{bmatrix} f(x_1) & \dots & f(x_i) \end{bmatrix}$$

$$= f(X)$$

### Backward Propagation

Étant donné  $\partial E/\partial Y$ , nous voulons calculer  $\partial E/\partial X$ .

$$\begin{aligned} \frac{\partial E}{\partial X} &= \begin{bmatrix} \frac{\partial E}{\partial x_1} & \dots & \frac{\partial E}{\partial x_i} \end{bmatrix} \\ &= \begin{bmatrix} \frac{\partial E}{\partial y_1} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial E}{\partial y_i} \frac{\partial y_i}{\partial x_i} \end{bmatrix} \\ &= \begin{bmatrix} \frac{\partial E}{\partial y_1} f'(x_1) & \dots & \frac{\partial E}{\partial y_i} f'(x_i) \end{bmatrix} \\ &= \begin{bmatrix} \frac{\partial E}{\partial y_1} & \dots & \frac{\partial E}{\partial y_i} \end{bmatrix} \odot [f'(x_1) \dots f'(x_i)] \\ &= \frac{\partial E}{\partial Y} \odot f'(X) \end{aligned}$$

Attention, nous utilisons ici une multiplication scalaire (élément par élément) entre les deux matrices.

### Coder la couche d'activation

Le code pour la couche d'activation est aussi très simple.

```

1 from layer import Layer
2
3 # inherit from base class Layer
4 class ActivationLayer(Layer):
5     def __init__(self, activation, activation_prime):
6         self.activation = activation
7         self.activation_prime = activation_prime
8
9     # returns the activated input
10    def forward_propagation(self, input_data):
11        self.input = input_data
12        self.output = self.activation(self.input)
13        return self.output
14
15    # Returns input_error=dE/dX for a given output_error=dE/dY.
16    # learning_rate is not used because there is no "learnable" parameters.
17    def backward_propagation(self, output_error, learning_rate):
18        return self.activation_prime(self.input) * output_error

```

Vous pouvez également écrire certaines fonctions d'activation et leurs dérivées dans un fichier séparé. Elles seront utilisées plus tard pour créer une couche d'activation.

```

1 import numpy as np
2
3 # activation function and its derivative
4 def tanh(x):
5     return np.tanh(x);
6
7 def tanh_prime(x):
8     return 1-np.tanh(x)**2;

```



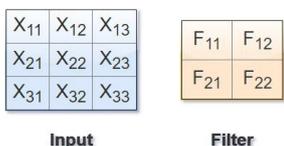
# Couche convolution

**Fondamental : Propagation avant et arrière dans un réseau neuronal convolutif**

Nous démontrons l'utilisation de l'opération de convolution pour effectuer la propagation arrière dans un CNN

*Différence entre l'opération de convolution et la corrélation.*

Considérons l'entrée et le filtre qui va être utilisé pour effectuer la convolution comme indiqué ci-dessous.

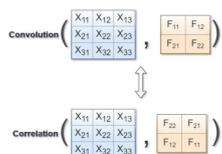


Ensuite, la corrélation de la matrice de filtre avec la matrice d'entrée est décrite dans la figure ci-dessous

$$\begin{pmatrix} O_{11} & O_{12} \\ O_{21} & O_{22} \end{pmatrix} = \text{Correlation} \left( \begin{pmatrix} X_{11} & X_{12} & X_{13} \\ X_{21} & X_{22} & X_{23} \\ X_{31} & X_{32} & X_{33} \end{pmatrix}, \begin{pmatrix} F_{11} & F_{12} \\ F_{21} & F_{22} \end{pmatrix} \right)$$

$$\begin{aligned}
 O_{11} &= F_{11}X_{11} + F_{12}X_{12} + F_{21}X_{21} + F_{22}X_{22} \\
 O_{12} &= F_{11}X_{12} + F_{12}X_{13} + F_{21}X_{22} + F_{22}X_{23} \\
 O_{21} &= F_{11}X_{21} + F_{12}X_{22} + F_{21}X_{31} + F_{22}X_{32} \\
 O_{22} &= F_{11}X_{22} + F_{12}X_{23} + F_{21}X_{32} + F_{22}X_{33}
 \end{aligned}$$

Maintenant, la convolution de la matrice de filtre avec l'image d'entrée est la même que la rotation du filtre de 180 degrés et ensuite l'exécution de la corrélation de la matrice de filtre tournée avec la matrice d'entrée.



Comme on peut le voir sur l'image ci-dessus, l'opération de convolution est la même que celle de l'opération de corrélation mais avec un filtre tourné.

*Propagation vers l'avant et vers l'arrière à l'aide de l'opération de convolution.*

Remarque: Pour dériver l'équation des gradients pour les valeurs de filtre et les valeurs de matrice d'entrée, nous considérerons que l'opération de convolution est la même que l'opération de corrélation, juste pour plus de simplicité.

Par conséquent, l'opération de convolution peut être écrite comme décrit dans la figure ci-dessous.

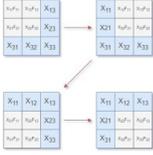


Couche convolution

$$\begin{pmatrix} O_{11} & O_{12} \\ O_{21} & O_{22} \end{pmatrix} = \text{Convolution} \left( \begin{pmatrix} X_{11} & X_{12} & X_{13} \\ X_{21} & X_{22} & X_{23} \\ X_{31} & X_{32} & X_{33} \end{pmatrix}, \begin{pmatrix} F_{11} & F_{12} \\ F_{21} & F_{22} \end{pmatrix} \right)$$

$$\begin{aligned} O_{11} &= F_{11}X_{11} + F_{12}X_{12} + F_{21}X_{21} + F_{22}X_{22} \\ O_{12} &= F_{11}X_{12} + F_{12}X_{13} + F_{21}X_{22} + F_{22}X_{23} \\ O_{21} &= F_{11}X_{21} + F_{12}X_{22} + F_{21}X_{31} + F_{22}X_{32} \\ O_{22} &= F_{11}X_{22} + F_{12}X_{23} + F_{21}X_{32} + F_{22}X_{33} \end{aligned}$$

Il peut être visualisé dans la figure ci-dessous.



Maintenant, pour calculer les gradients du filtre «F» par rapport à l'erreur «E», les équations suivantes doivent être résolues.

$$\begin{aligned} \frac{\partial E}{\partial F_{11}} &= \frac{\partial E}{\partial O_{11}} \frac{\partial O_{11}}{\partial F_{11}} + \frac{\partial E}{\partial O_{12}} \frac{\partial O_{12}}{\partial F_{11}} + \frac{\partial E}{\partial O_{21}} \frac{\partial O_{21}}{\partial F_{11}} + \frac{\partial E}{\partial O_{22}} \frac{\partial O_{22}}{\partial F_{11}} \\ \frac{\partial E}{\partial F_{12}} &= \frac{\partial E}{\partial O_{11}} \frac{\partial O_{11}}{\partial F_{12}} + \frac{\partial E}{\partial O_{12}} \frac{\partial O_{12}}{\partial F_{12}} + \frac{\partial E}{\partial O_{21}} \frac{\partial O_{21}}{\partial F_{12}} + \frac{\partial E}{\partial O_{22}} \frac{\partial O_{22}}{\partial F_{12}} \\ \frac{\partial E}{\partial F_{21}} &= \frac{\partial E}{\partial O_{11}} \frac{\partial O_{11}}{\partial F_{21}} + \frac{\partial E}{\partial O_{12}} \frac{\partial O_{12}}{\partial F_{21}} + \frac{\partial E}{\partial O_{21}} \frac{\partial O_{21}}{\partial F_{21}} + \frac{\partial E}{\partial O_{22}} \frac{\partial O_{22}}{\partial F_{21}} \\ \frac{\partial E}{\partial F_{22}} &= \frac{\partial E}{\partial O_{11}} \frac{\partial O_{11}}{\partial F_{22}} + \frac{\partial E}{\partial O_{12}} \frac{\partial O_{12}}{\partial F_{22}} + \frac{\partial E}{\partial O_{21}} \frac{\partial O_{21}}{\partial F_{22}} + \frac{\partial E}{\partial O_{22}} \frac{\partial O_{22}}{\partial F_{22}} \end{aligned}$$

qui évalue à

$$\begin{aligned} \frac{\partial E}{\partial F_{11}} &= \frac{\partial E}{\partial O_{11}} X_{11} + \frac{\partial E}{\partial O_{12}} X_{12} + \frac{\partial E}{\partial O_{21}} X_{21} + \frac{\partial E}{\partial O_{22}} X_{22} \\ \frac{\partial E}{\partial F_{12}} &= \frac{\partial E}{\partial O_{11}} X_{12} + \frac{\partial E}{\partial O_{12}} X_{13} + \frac{\partial E}{\partial O_{21}} X_{22} + \frac{\partial E}{\partial O_{22}} X_{23} \\ \frac{\partial E}{\partial F_{21}} &= \frac{\partial E}{\partial O_{11}} X_{21} + \frac{\partial E}{\partial O_{12}} X_{22} + \frac{\partial E}{\partial O_{21}} X_{31} + \frac{\partial E}{\partial O_{22}} X_{32} \\ \frac{\partial E}{\partial F_{22}} &= \frac{\partial E}{\partial O_{11}} X_{22} + \frac{\partial E}{\partial O_{12}} X_{23} + \frac{\partial E}{\partial O_{21}} X_{32} + \frac{\partial E}{\partial O_{22}} X_{33} \end{aligned}$$

Si nous regardons de près cette équation ci-dessus peut être écrite sous la forme de notre opération de convolution.

$$\begin{pmatrix} \frac{\partial E}{\partial F_{11}} & \frac{\partial E}{\partial F_{12}} \\ \frac{\partial E}{\partial F_{21}} & \frac{\partial E}{\partial F_{22}} \end{pmatrix} = \text{Convolution} \left( \begin{pmatrix} X_{11} & X_{12} & X_{13} \\ X_{21} & X_{22} & X_{23} \\ X_{31} & X_{32} & X_{33} \end{pmatrix}, \begin{pmatrix} \frac{\partial E}{\partial O_{11}} & \frac{\partial E}{\partial O_{12}} \\ \frac{\partial E}{\partial O_{21}} & \frac{\partial E}{\partial O_{22}} \end{pmatrix} \right)$$

De même, nous pouvons trouver les gradients de la matrice d'entrée «X» par rapport à l'erreur «E».

$$\begin{aligned} \frac{\partial E}{\partial X_{11}} &= \frac{\partial E}{\partial O_{11}} F_{11} + \frac{\partial E}{\partial O_{12}} F_{12} + \frac{\partial E}{\partial O_{21}} F_{21} + \frac{\partial E}{\partial O_{22}} F_{22} \\ \frac{\partial E}{\partial X_{12}} &= \frac{\partial E}{\partial O_{11}} F_{12} + \frac{\partial E}{\partial O_{12}} F_{13} + \frac{\partial E}{\partial O_{21}} F_{22} + \frac{\partial E}{\partial O_{22}} F_{23} \\ \frac{\partial E}{\partial X_{13}} &= \frac{\partial E}{\partial O_{11}} F_{13} + \frac{\partial E}{\partial O_{12}} F_{23} + \frac{\partial E}{\partial O_{21}} F_{22} + \frac{\partial E}{\partial O_{22}} F_{23} \\ \frac{\partial E}{\partial X_{21}} &= \frac{\partial E}{\partial O_{11}} F_{21} + \frac{\partial E}{\partial O_{12}} F_{22} + \frac{\partial E}{\partial O_{21}} F_{31} + \frac{\partial E}{\partial O_{22}} F_{32} \\ \frac{\partial E}{\partial X_{22}} &= \frac{\partial E}{\partial O_{11}} F_{22} + \frac{\partial E}{\partial O_{12}} F_{23} + \frac{\partial E}{\partial O_{21}} F_{32} + \frac{\partial E}{\partial O_{22}} F_{33} \\ \frac{\partial E}{\partial X_{23}} &= \frac{\partial E}{\partial O_{11}} F_{23} + \frac{\partial E}{\partial O_{12}} F_{33} + \frac{\partial E}{\partial O_{21}} F_{32} + \frac{\partial E}{\partial O_{22}} F_{33} \\ \frac{\partial E}{\partial X_{31}} &= \frac{\partial E}{\partial O_{21}} F_{21} + \frac{\partial E}{\partial O_{22}} F_{22} \\ \frac{\partial E}{\partial X_{32}} &= \frac{\partial E}{\partial O_{21}} F_{22} + \frac{\partial E}{\partial O_{22}} F_{23} \\ \frac{\partial E}{\partial X_{33}} &= \frac{\partial E}{\partial O_{21}} F_{23} + \frac{\partial E}{\partial O_{22}} F_{33} \end{aligned}$$

Maintenant, le calcul ci-dessus peut être obtenu par un autre type d'opération de convolution connue sous le nom de convolution complète.

Afin d'obtenir les gradients de la matrice d'entrée, nous devons faire pivoter le filtre de 180 degrés et calculer la convolution complète du filtre tourné par les gradients de la sortie par rapport à l'erreur, comme représenté dans l'image ci-dessous.

$$\begin{pmatrix} \partial E / \partial X_{11} & \partial E / \partial X_{12} & \partial E / \partial X_{13} \\ \partial E / \partial X_{21} & \partial E / \partial X_{22} & \partial E / \partial X_{23} \\ \partial E / \partial X_{31} & \partial E / \partial X_{32} & \partial E / \partial X_{33} \end{pmatrix} = \text{Full\_Convolution} \left( \begin{pmatrix} \partial E / \partial O_{11} & \partial E / \partial O_{12} \\ \partial E / \partial O_{21} & \partial E / \partial O_{22} \end{pmatrix}, \begin{pmatrix} F_{22} & F_{21} \\ F_{12} & F_{11} \end{pmatrix} \right)$$

La convolution complète peut être visualisée comme exécutant la procédure comme représenté dans la figure ci-dessous.



- Par conséquent, à la fois la propagation vers l'avant et vers l'arrière peut être effectuée en utilisant l'opération de convolution.
- Pour calculer les gradients des couches de regroupement et de Relu, les gradients peuvent être calculés en suivant la même procédure d'utilisation de la règle de chaîne des dérivés.

la couche de convolution

```

1 from layer import Layer
2 from scipy import signal
3 import numpy as np
4
5 ## Math behind this layer can found at :
6 ## https://medium.com/@2017csm1006/forward-and-backpropagation-in-convolutional-
7   neural-network-4dfa96d7b37e
8 # inherit from base class Layer
9 # This convolutional layer is always with stride 1
10 class ConvLayer(Layer):
11     # input_shape = (i,j,d)
12     # kernel_shape = (m,n)
13     # layer_depth = output_depth
14     def __init__(self, input_shape, kernel_shape, layer_depth):
15         self.input_shape = input_shape
16         self.input_depth = input_shape[2]
17         self.kernel_shape = kernel_shape
18         self.layer_depth = layer_depth
19         self.output_shape = (input_shape[0]-kernel_shape[0]+1, input_shape[1]-
20 kernel_shape[1]+1, layer_depth)
21         self.weights = np.random.rand(kernel_shape[0], kernel_shape[1], self.
22 input_depth, layer_depth) - 0.5
23         self.bias = np.random.rand(layer_depth) - 0.5
24
25     # returns output for a given input
26     def forward_propagation(self, input):
27         self.input = input
28         self.output = np.zeros(self.output_shape)
29
30         for k in range(self.layer_depth):
31             for d in range(self.input_depth):
32                 self.output[:, :, k] += signal.correlate2d(self.input[:, :, d], self.
33 weights[:, :, d, k], 'valid') + self.bias[k]
34
35         return self.output

```

```

34 # computes dE/dW, dE/dB for a given output_error=dE/dY. Returns input_error=dE
    /dx.
35 def backward_propagation(self, output_error, learning_rate):
36     in_error = np.zeros(self.input_shape)
37     dWeights = np.zeros((self.kernel_shape[0], self.kernel_shape[1], self.
input_depth, self.layer_depth))
38     dBias = np.zeros(self.layer_depth)
39
40     for k in range(self.layer_depth):
41         for d in range(self.input_depth):
42             in_error[:, :, d] += signal.convolve2d(output_error[:, :, k], self.
weights[:, :, d, k], 'full')
43             dWeights[:, :, d, k] = signal.correlate2d(self.input[:, :, d],
output_error[:, :, k], 'valid')
44             dBias[k] = self.layer_depth * np.sum(output_error[:, :, k])
45
46     self.weights -= learning_rate*dWeights
47     self.bias -= learning_rate*dBias
48     return in_error

```

### exemple

```

1 import numpy as np
2
3 from network import Network
4 from conv_layer import ConvLayer
5 from activation_layer import ActivationLayer
6 from activations import tanh, tanh_prime
7 from losses import mse, mse_prime
8
9 # training data
10 x_train = [np.random.rand(10,10,1)]
11 y_train = [np.random.rand(4,4,2)]
12
13 # network
14 net = Network()
15 net.add(ConvLayer((10,10,1), (3,3), 1))
16 net.add(ActivationLayer(tanh, tanh_prime))
17 net.add(ConvLayer((8,8,1), (3,3), 1))
18 net.add(ActivationLayer(tanh, tanh_prime))
19 net.add(ConvLayer((6,6,1), (3,3), 2))
20 net.add(ActivationLayer(tanh, tanh_prime))
21
22 # train
23 net.use(mse, mse_prime)
24 net.fit(x_train, y_train, epochs=1000, learning_rate=0.3)
25
26 # test
27 out = net.predict(x_train)
28 print("predicted = ", out)
29 print("expected = ", y_train)

```

# Flatten layer

VIII

- Les couches de convolution 2D traitant des données 2D (par exemple, des images) produisent généralement un tenseur tridimensionnel, les dimensions étant la résolution de l'image (moins la taille du filtre -1) et le nombre de filtres.
- Si vous utilisez un calque avec  $N$  filtres de taille  $s$  sur des images  $L \times H$ , le résultat sera de dimension  $(L - (s - 1), H - (s - 1), N)$ .
- Vous avez  $N$  images de résolution  $(L - (s - 1), H - (s - 1))$ .
- Cette structure est importante si vous souhaitez enchaîner des couches de convolution entre elles ou avec d'autres couches qui effectuent un traitement spatial (pooling, upscaling, etc.).
- Mais si vous voulez faire de la classification, dans les dernières étapes de votre réseau, vous voulez généralement utiliser des couches entièrement connectées qui ne prennent en compte aucune structure (spatiale ou autre) pour le traitement, vous voulez donc juste la sortie de dernières couches de convolution à être considérées comme un gros morceau de données non structurées.
- C'est ce que signifie «aplatir». Il brise la structure spatiale des données et transforme votre tenseur tridimensionnel  $(L - (s - 1), H - (s - 1), N)$  en un tenseur monodimensionnel (un vecteur) de taille  $(L - (s - 1)) \times (H - (s - 1)) \times N$ .

```

1 from layer import Layer
2
3 # inherit from base class Layer
4 class FlattenLayer(Layer):
5     # returns the flattened input
6     def forward_propagation(self, input_data):
7         self.input = input_data
8         self.output = input_data.flatten().reshape((1,-1))
9         return self.output
10
11     # Returns input_error=dE/dX for a given output_error=dE/dY.
12     # learning_rate is not used because there is no "learnable" parameters.
13     def backward_propagation(self, output_error, learning_rate):
14         return output_error.reshape(self.input.shape)

```

```

1 import numpy as np
2
3 from network import Network
4 from fc_layer import FCLayer
5 from conv_layer import ConvLayer
6 from flatten_layer import FlattenLayer
7 from activation_layer import ActivationLayer
8 from activations import tanh, tanh_prime
9 from losses import mse, mse_prime
10
11 from keras.datasets import mnist
12 from keras.utils import np_utils
13

```

```

14 # load MNIST from server
15 (x_train, y_train), (x_test, y_test) = mnist.load_data()
16
17 # training data : 60000 samples
18 # reshape and normalize input data
19 x_train = x_train.reshape(x_train.shape[0], 28, 28, 1)
20 x_train = x_train.astype('float32')
21 x_train /= 255
22 # encode output which is a number in range [0,9] into a vector of size 10
23 # e.g. number 3 will become [0, 0, 0, 1, 0, 0, 0, 0, 0, 0]
24 y_train = np_utils.to_categorical(y_train)
25
26 # same for test data : 10000 samples
27 x_test = x_test.reshape(x_test.shape[0], 28, 28, 1)
28 x_test = x_test.astype('float32')
29 x_test /= 255
30 y_test = np_utils.to_categorical(y_test)
31
32 # Network
33 net = Network()
34 net.add(ConvLayer((28, 28, 1), (3, 3), 1)) # input_shape=(28, 28, 1) ;
      output_shape=(26, 26, 1)
35 net.add(ActivationLayer(tanh, tanh_prime))
36 net.add(FlattenLayer()) # input_shape=(26, 26, 1) ;
      output_shape=(1, 26*26*1)
37 net.add(FCLayer(26*26*1, 100)) # input_shape=(1, 26*26*1) ;
      output_shape=(1, 100)
38 net.add(ActivationLayer(tanh, tanh_prime))
39 net.add(FCLayer(100, 10)) # input_shape=(1, 100) ;
      output_shape=(1, 10)
40 net.add(ActivationLayer(tanh, tanh_prime))
41
42 # train on 1000 samples
43 # as we didn't implemented mini-batch GD, training will be pretty slow if we
      update at each iteration on 60000 samples...
44 net.use(mse, mse_prime)
45 net.fit(x_train[0:1000], y_train[0:1000], epochs=100, learning_rate=0.1)
46
47 # test on 3 samples
48 out = net.predict(x_test[0:3])
49 print("\n")
50 print("predicted values : ")
51 print(out, end="\n")
52 print("true values : ")
53 print(y_test[0:3])

```

# Fonction de perte

IX

## Fonction de perte

Jusqu'à présent, pour une couche donnée, nous supposons que  $\partial E/\partial Y$  était donné (par la couche suivante). Mais qu'advient-il de la dernière couche? Comment obtient-elle  $\partial E/\partial Y$ ? Nous le donnons simplement manuellement, et cela dépend de la façon dont nous définissons l'erreur.

L'erreur du réseau, qui mesure le degré de performance du modèle pour une entrée donnée, est définie par nous-même. Il existe de nombreuses façons de définir l'erreur, et l'une des plus connues est appelée MSE — Mean Squared Error.

$$E = \frac{1}{n} \sum_i^n (y_i^* - y_i)^2$$

Où  $y^*$  et  $y$  désignent respectivement la sortie souhaitée et la sortie obtenue. Vous pouvez penser à la perte comme une dernière couche qui regroupe tous les neurones de sortie et les écrase en un seul neurone. Ce dont nous avons besoin maintenant, comme pour toutes les autres couches, c'est de définir  $\partial E/\partial Y$ . Excepté que maintenant, nous avons enfin "atteint" E!

$$\begin{aligned} \frac{\partial E}{\partial Y} &= \left[ \frac{\partial E}{\partial y_1} \quad \dots \quad \frac{\partial E}{\partial y_i} \right] \\ &= \frac{2}{n} [y_1 - y_1^* \quad \dots \quad y_i - y_i^*] \\ &= \frac{2}{n} (Y - Y^*) \end{aligned}$$

Et voilà! Il suffira de donner cette valeur à la dernière couche lors de la passe arrière, ce qui lui permettra d'ajuster ses paramètres, puis elle calculera  $\partial E/\partial X$  qu'elle passera à la couche d'avant, qui fera le même procédé à son tour, etc.

Nous pouvons implémenter la fonction de perte dans un fichier séparé. Elle sera utilisée lors de la création du réseau.

```
1 import numpy as np
2
3 # loss function and its derivative
4 def mse(y_true, y_pred):
5     return np.mean(np.power(y_true-y_pred, 2));
6
7 def mse_prime(y_true, y_pred):
8     return 2*(y_pred-y_true)/y_true.size;
```

# Classe Network



## Classe Network

Bientôt fini, tenez bon ! Maintenant que nous avons tous nos blocs de code prêts à l'emploi, nous allons faire une classe appelée Network qui permettra de construire ces fameux réseaux de neurones !

J'ai commenté presque chaque partie du code, il ne devrait pas être trop compliqué à comprendre si vous avez saisi les étapes précédentes. Néanmoins, laissez un commentaire si vous avez des questions, je répondrai avec plaisir !

```

1 class Network:
2     def __init__(self):
3         self.layers = []
4         self.loss = None
5         self.loss_prime = None
6
7     # add layer to network
8     def add(self, layer):
9         self.layers.append(layer)
10
11    # set loss to use
12    def use(self, loss, loss_prime):
13        self.loss = loss
14        self.loss_prime = loss_prime
15
16    # predict output for given input
17    def predict(self, input_data):
18        # sample dimension first
19        samples = len(input_data)
20        result = []
21
22        # run network over all samples
23        for i in range(samples):
24            # forward propagation
25            output = input_data[i]
26            for layer in self.layers:
27                output = layer.forward_propagation(output)
28            result.append(output)
29
30        return result
31
32    # train the network
33    def fit(self, x_train, y_train, epochs, learning_rate):
34        # sample dimension first
35        samples = len(x_train)
36
37        # training loop

```

```
38     for i in range(epochs):
39         err = 0
40         for j in range(samples):
41             # forward propagation
42             output = x_train[j]
43             for layer in self.layers:
44                 output = layer.forward_propagation(output)
45
46             # compute loss (for display purpose only)
47             err += self.loss(y_train[j], output)
48
49             # backward propagation
50             error = self.loss_prime(y_train[j], output)
51             for layer in reversed(self.layers):
52                 error = layer.backward_propagation(error, learning_rate)
53
54             # calculate average error on all samples
55             err /= samples
56             print('epoch %d/%d  error=%f' % (i+1, epochs, err))
```

# Construire un réseau de neurone


 XI

## Construire un réseau de neurone

Construire un réseau de neurone

Enfin ! Nous pouvons utiliser notre classe pour créer un réseau de neurones avec autant de couches que nous voulons ! Nous allons construire deux réseaux de neurones: un simple XOR et un solveur MNIST.

### Exemple : Résoudre XOR

Commencer par un XOR est toujours important car c'est un moyen simple de savoir si le réseau apprend quelque chose.

```

1 import numpy as np
2
3 from network import Network
4 from fc_layer import FCLayer
5 from activation_layer import ActivationLayer
6 from activations import tanh, tanh_prime
7 from losses import mse, mse_prime
8
9 # training data
10 x_train = np.array([[0,0], [0,1], [1,0], [1,1]])
11 y_train = np.array([[0], [1], [1], [0]])
12
13 # network
14 net = Network()
15 net.add(FCLayer(2, 3))
16 net.add(ActivationLayer(tanh, tanh_prime))
17 net.add(FCLayer(3, 1))
18 net.add(ActivationLayer(tanh, tanh_prime))
19
20 # train
21 net.use(mse, mse_prime)
22 net.fit(x_train, y_train, epochs=1000, learning_rate=0.1)
23
24 # test
25 out = net.predict(x_train)
26 print(out)

```

```

1 $ python xor.py
2 epoch 1/1000 error=0.322980
3 epoch 2/1000 error=0.311174
4 epoch 3/1000 error=0.307195
5 ...

```

```
6 epoch 998/1000 error=0.000243
7 epoch 999/1000 error=0.000242
8 epoch 1000/1000 error=0.000242[
9   array([[ 0.00077435]]),
10  array([[ 0.97760742]]),
11  array([[ 0.97847793]]),
12  array([[ -0.00131305]])
13 ]
```

### Remarque

---

Je ne pense pas avoir besoin d'insister sur beaucoup de choses. Faites attention avec les données d'entraînement, vous devriez toujours avoir la dimension de l'échantillon en premier. Par exemple, avec le problème xor, la dimension des données d'entrées devrait être (4,1,2).

# exemple MNIST


 XII

ddd

```

1 import numpy as np
2
3 from network import Network
4 from fc_layer import FCLayer
5 from activation_layer import ActivationLayer
6 from activations import tanh, tanh_prime
7 from losses import mse, mse_prime
8
9 from keras.datasets import mnist
10 from keras.utils import np_utils
11
12 # load MNIST from server
13 (x_train, y_train), (x_test, y_test) = mnist.load_data()
14
15 # training data : 60000 samples
16 # reshape and normalize input data
17 x_train = x_train.reshape(x_train.shape[0], 1, 28*28)
18 x_train = x_train.astype('float32')
19 x_train /= 255
20 # encode output which is a number in range [0,9] into a vector of size 10
21 # e.g. number 3 will become [0, 0, 0, 1, 0, 0, 0, 0, 0, 0]
22 y_train = np_utils.to_categorical(y_train)
23
24 # same for test data : 10000 samples
25 x_test = x_test.reshape(x_test.shape[0], 1, 28*28)
26 x_test = x_test.astype('float32')
27 x_test /= 255
28 y_test = np_utils.to_categorical(y_test)
29
30 # Network
31 net = Network()
32 net.add(FCLayer(28*28, 100))           # input_shape=(1, 28*28) ;
33   output_shape=(1, 100)
34 net.add(ActivationLayer(tanh, tanh_prime))
35 net.add(FCLayer(100, 50))             # input_shape=(1, 100) ;
36   output_shape=(1, 50)
37 net.add(ActivationLayer(tanh, tanh_prime))
38 net.add(FCLayer(50, 10))              # input_shape=(1, 50) ;
39   output_shape=(1, 10)
40 net.add(ActivationLayer(tanh, tanh_prime))
41
42 # train on 1000 samples
43 # as we didn't implemented mini-batch GD, training will be pretty slow if we
44   update at each iteration on 60000 samples...
45 net.use(mse, mse_prime)
46 net.fit(x_train[0:1000], y_train[0:1000], epochs=35, learning_rate=0.1)

```

```
43
44 # test on 3 samples
45 out = net.predict(x_test[0:3])
46 print("\n")
47 print("predicted values : ")
48 print(out, end="\n")
49 print("true values : ")
50 print(y_test[0:3])
51
```

# Le deep learning

XIII

cc

```

1 import numpy as np
2
3 from network import Network
4 from fc_layer import FCLayer
5 from conv_layer import ConvLayer
6 from flatten_layer import FlattenLayer
7 from activation_layer import ActivationLayer
8 from activations import tanh, tanh_prime
9 from losses import mse, mse_prime
10
11 from keras.datasets import mnist
12 from keras.utils import np_utils
13
14 # load MNIST from server
15 (x_train, y_train), (x_test, y_test) = mnist.load_data()
16
17 # training data : 60000 samples
18 # reshape and normalize input data
19 x_train = x_train.reshape(x_train.shape[0], 28, 28, 1)
20 x_train = x_train.astype('float32')
21 x_train /= 255
22 # encode output which is a number in range [0,9] into a vector of size 10
23 # e.g. number 3 will become [0, 0, 0, 1, 0, 0, 0, 0, 0, 0]
24 y_train = np_utils.to_categorical(y_train)
25
26 # same for test data : 10000 samples
27 x_test = x_test.reshape(x_test.shape[0], 28, 28, 1)
28 x_test = x_test.astype('float32')
29 x_test /= 255
30 y_test = np_utils.to_categorical(y_test)
31
32 # Network
33 net = Network()
34 net.add(ConvLayer((28, 28, 1), (3, 3), 1)) # input_shape=(28, 28, 1) ;
35     output_shape=(26, 26, 1)
36 net.add(ActivationLayer(tanh, tanh_prime))
37 net.add(FlattenLayer()) # input_shape=(26, 26, 1) ;
38     output_shape=(1, 26*26*1)
39 net.add(FCLayer(26*26*1, 100)) # input_shape=(1, 26*26*1) ;
40     output_shape=(1, 100)
41 net.add(ActivationLayer(tanh, tanh_prime))
42 net.add(FCLayer(100, 10)) # input_shape=(1, 100) ;
43     output_shape=(1, 10)
44 net.add(ActivationLayer(tanh, tanh_prime))
45
46 # train on 1000 samples

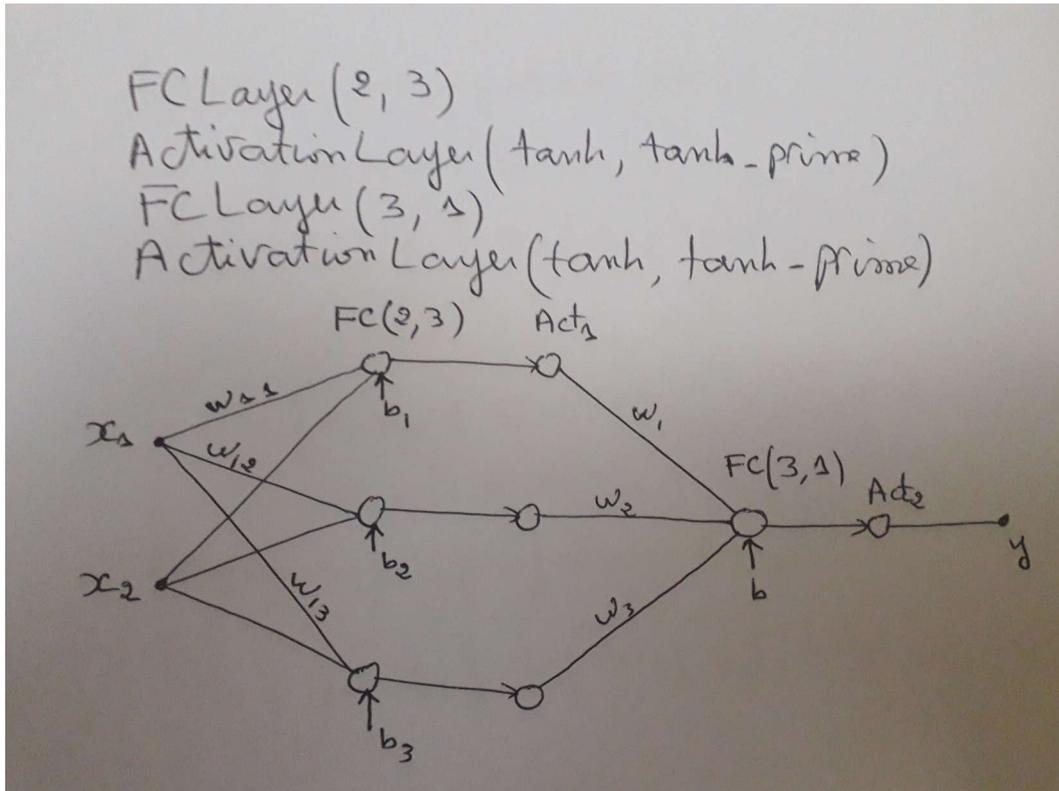
```

```
43 # as we didn't implemented mini-batch GD, training will be pretty slow if we
    update at each iteration on 60000 samples...
44 net.use(mse, mse_prime)
45 net.fit(x_train[0:1000], y_train[0:1000], epochs=100, learning_rate=0.1)
46
47 # test on 3 samples
48 out = net.predict(x_test[0:3])
49 print("\n")
50 print("predicted values : ")
51 print(out, end="\n")
52 print("true values : ")
53 print(y_test[0:3])
```

# Exercice :

XIV

Dessiner le schéma de réseau de neurones de l'exemple XOR et explique comment l'apprentissage et la prédiction se font



forward propagation FC(3,1):

$$y = xW + B$$

$$y = [x_1, x_2, x_3] \times \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} + [b] =$$

$$x_1 w_1 + x_2 w_2 + x_3 w_3 + b$$

forward propagation Act<sub>2</sub>:

$$y = \text{Act}_2(x_1 w_1 + x_2 w_2 + x_3 w_3 + b)$$

forward propagation FC(2,3):

$$Y = XW + B$$

$$[y_1, y_2, y_3] = [x_1, x_2] \cdot \begin{bmatrix} w_{11}, w_{12}, w_{13} \\ w_{21}, w_{22}, w_{23} \end{bmatrix} +$$

$$[b_1, b_2, b_3]$$

$$= [x_1 w_{11} + x_2 w_{21} + b_1, x_1 w_{12} + x_2 w_{22} + b_2 +$$

$$x_1 w_{13} + x_2 w_{23} + b_3]$$

forward propagation Act<sub>1</sub>:

$$[y_1, y_2, y_3] = [\text{Act}_1(y_1), \text{Act}_1(y_2), \text{Act}_1(y_3)]$$

backward propagation  $FC(3,1)$

$$\frac{\partial E}{\partial b} = error = e * (y - y_{true}) = \frac{\partial E}{\partial y}$$

$$\frac{\partial E}{\partial x} = input\_error = [error] * \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}$$

$$= [error * w_1 + error * w_2 + error * w_3]$$

$$= \frac{\partial E}{\partial x} = \frac{\partial E}{\partial y} * W^T$$

$$\frac{\partial E}{\partial w} = X^t * \frac{\partial E}{\partial y} = weight\_error = \begin{bmatrix} x_1 * error \\ x_2 * error \\ x_3 * error \end{bmatrix}$$

update parameters:

weights = weight - learning\\_rate \* weight\\_error

$$= [w_1, w_2, w_3] - learning\_rate * \begin{bmatrix} x_1 * error \\ x_2 * error \\ x_3 * error \end{bmatrix}$$

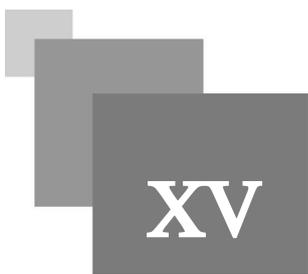
$$= [w_1 - \alpha * x_1 * error, w_2 - \alpha * x_2 * error, w_3 - \alpha * x_3 * error]$$

bias =  $\alpha * output\_error$

$$= \alpha * e * (y - y_{true})$$

Return input\\_error  $\left(\frac{\partial E}{\partial x}\right) \Rightarrow error$

# Exercice :



## Question

Les couches de convolution 2D traitant des données 2D (par exemple, des images) produisent généralement un tenseur tridimensionnel, les dimensions étant la résolution de l'image (moins la taille du filtre -1) et le nombre de filtres.

Si vous utilisez un calque avec  $N$  filtres de taille  $s$  sur des images  $L \times H$ , le résultat sera de dimension  $(L - (s - 1), H - (s - 1), N)$ .

Vous avez  $N$  images de résolution  $(L - (s - 1), H - (s - 1))$ .

Cette structure est importante si vous souhaitez enchaîner des couches de convolution entre elles ou avec d'autres couches qui effectuent un traitement spatial (pooling, upscaling, etc.).

Si  $L = H = 4$ ,  $N = 2$  alors nous avons 2 images de résolution  $3 \times 3$

La taille du filtre est  $s = 2$

$$x_{train}[0](3, 3, 2) = \left[ \begin{bmatrix} 1 & 1.5 \\ 2 & 2.5 \\ 3 & 3.5 \end{bmatrix}, \begin{bmatrix} 4 & 4.5 \\ 5 & 5.5 \\ 6 & 6.5 \end{bmatrix}, \begin{bmatrix} 7 & 7.5 \\ 8 & 8.5 \\ 9 & 9.6 \end{bmatrix} \right]$$

$$y_{train}(2, 2) = \left[ \begin{bmatrix} 90 \\ 112 \end{bmatrix}, \begin{bmatrix} 156 \\ 177 \end{bmatrix} \right]$$

Comme nous avons deux images, il y a deux filtres. Le filtre est  $2 \times 2$

$$weights = \left[ \left[ \begin{bmatrix} 1 \\ 1.5 \end{bmatrix}, \begin{bmatrix} 2 \\ 2.5 \end{bmatrix} \right], \left[ \begin{bmatrix} 3 \\ 3.5 \end{bmatrix}, \begin{bmatrix} 4 \\ 4.5 \end{bmatrix} \right] \right]$$

Nous avons une seule sortie (profondeur de la couche = 1) :

$$B = [2]$$

Forme de output =  $(3-2+1, 3-2+1, 1) = (2, 2, 1)$

$$x_{train}[0][:, :, 0] = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}$$

$$weights[:, :, 0, 0] = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$$

$$output1 = correlation\left(\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}, \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}\right) + B = \begin{bmatrix} 39 & 49 \\ 69 & 79 \end{bmatrix}$$

$$x_{train}[0][:, :, 1] = \begin{bmatrix} 1.5 & 2.5 & 3.5 \\ 4.5 & 5.5 & 6.5 \\ 7.5 & 8.5 & 9.5 \end{bmatrix}$$

$$weights[:, :, 1, 0] = \begin{bmatrix} 1.5 & 2.5 \\ 3.5 & 4.5 \end{bmatrix}$$

$$output2 = correlation\left(\begin{bmatrix} 1.5 & 2.5 & 3.5 \\ 4.5 & 5.5 & 6.5 \\ 7.5 & 8.5 & 9.5 \end{bmatrix}, \begin{bmatrix} 1.5 & 2.5 \\ 3.5 & 4.5 \end{bmatrix}\right) + B = \begin{bmatrix} 51 & 63 \\ 87 & 99 \end{bmatrix}$$

$$output = output1 + output2 = \begin{bmatrix} 90 & 112 \\ 156 & 178 \end{bmatrix}$$

**Calcul de l'erreur E du réseau (MSE) :**

$$E = \frac{1}{n} \sum_{i=1}^n (y_i^* - y_i)^2 = \frac{1}{4} \left( \begin{bmatrix} 90 & 112 \\ 156 & 177 \end{bmatrix} - \begin{bmatrix} 90 & 112 \\ 156 & 178 \end{bmatrix} \right)^2 = 0.25$$

calcul des erreurs de sortie (dérivée de MSE) :

$$\frac{\partial E}{\partial Y} = \frac{2}{n} (Y - Y^*) = \frac{2}{4} \left( \begin{bmatrix} 90 & 112 \\ 156 & 178 \end{bmatrix} - \begin{bmatrix} 90 & 112 \\ 156 & 177 \end{bmatrix} \right) = 0.5$$

$$\frac{\partial E}{\partial W}[:, :, 0, 0] = correlation\left(\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}, \begin{bmatrix} 0 & 0 \\ 0 & 0.5 \end{bmatrix}\right) = \begin{bmatrix} 2.5 & 3 \\ 4 & 4.5 \end{bmatrix}$$

$$\frac{\partial E}{\partial W}[:, :, 1, 0] = correlation\left(\begin{bmatrix} 1.5 & 2.5 & 3.5 \\ 4.5 & 5.5 & 6.5 \\ 7.5 & 8.5 & 9.5 \end{bmatrix}, \begin{bmatrix} 0 & 0 \\ 0 & 0.5 \end{bmatrix}\right) = \begin{bmatrix} 2.75 & 3.25 \\ 4.25 & 4.75 \end{bmatrix}$$

$$\frac{\partial E}{\partial W} = \left[ \left[ \begin{bmatrix} 2.5 \\ 2.75 \end{bmatrix}, \begin{bmatrix} 3 \\ 3.25 \end{bmatrix} \right], \left[ \begin{bmatrix} 4 \\ 4.25 \end{bmatrix}, \begin{bmatrix} 4.5 \\ 4.75 \end{bmatrix} \right] \right]$$

$$\frac{\partial E}{\partial X}[:, :, 0] = convolution\left(\begin{bmatrix} 0 & 0 \\ 0 & 0.5 \end{bmatrix}, \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}\right) = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0.5 & 1 \\ 0 & 1.5 & 2 \end{bmatrix}$$

$$\frac{\partial E}{\partial X}[:, :, 1] = convolution\left(\begin{bmatrix} 0 & 0 \\ 0 & 0.5 \end{bmatrix}, \begin{bmatrix} 1.5 & 2.5 \\ 3.5 & 4.5 \end{bmatrix}\right) = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0.75 & 1.25 \\ 0 & 1.75 & 2.25 \end{bmatrix}$$

$$\frac{\partial E}{\partial X} = \left[ \left[ \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}, \begin{bmatrix} 0 & 0 \\ 0.5 & 0.75 \\ 1 & 1.25 \end{bmatrix}, \begin{bmatrix} 0 & 0 \\ 1.5 & 1.75 \\ 2 & 2.25 \end{bmatrix} \right] \right]$$

Mise à jour pour  $\alpha = 0.1$  :

$$weights = weights - \alpha * \frac{\partial E}{\partial W} = \left[ \left[ \begin{bmatrix} 1 \\ 1.5 \end{bmatrix}, \begin{bmatrix} 2 \\ 2.5 \end{bmatrix}, \begin{bmatrix} 3 \\ 3.5 \end{bmatrix}, \begin{bmatrix} 4 \\ 4.5 \end{bmatrix} \right] - 0.1 \times \left[ \left[ \begin{bmatrix} 2.5 \\ 2.75 \end{bmatrix}, \begin{bmatrix} 3 \\ 3.25 \end{bmatrix}, \begin{bmatrix} 4 \\ 4.25 \end{bmatrix}, \begin{bmatrix} 4.5 \\ 4.75 \end{bmatrix} \right] \right] =$$

$$bias = bias - \alpha * \frac{\partial E}{\partial B} = 2 - 0.1 \times 0.5 =$$

# Conclusion



Ce chapitre vous a montré comment programmer un réseau de neurones en Python pour réaliser un apprentissage automatique.

