

```
In [1]: import numpy as np
import matplotlib.pyplot as plot
```

# Fonctions d'activation

## Mise en situation

Toujours salarié à la poste, préposé aux algorithmes de reconnaissance des codes postaux :)

Le but du jour sera de faire évoluer un peu mon modèle pour voir si on peut faire mieux.

## Les fonctions d'activation

### Rappel

Il est nécessaire on l'a vu d'avoir dans nos réseaux de neurones une fonction d'activation pour casser la linéarité du calcul. On a utilisée jusque là la sigmoïde  $\sigma(x) = \frac{1}{1+e^{-x}}$

On va en introduire d'autres.

### Pourquoi la sigmoïde

La sigmoïde avait été introduite avec les régressions logistiques. Il nous fallait une fonction qui renvoie une probabilité à partir d'un nombre réel.

On aurait pu penser à un échelon Heaviside, mais c'est trop "plat". Avec une sigmoïde, on sait au moins de quel côté évoluer pour améliorer. Un échelon a une dérivée nulle presque partout et donc l'utilisation du gradient sera compromise.

### Tangente hyperbolique

Pour les cas où on n'a pas de besoin particulier d'une valeur entre 0 et 1, on peut utiliser la tangente hyperbolique. C'est une version un peu étirée de la sigmoïde, à valeur entre -1 et 1.

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

### ReLU

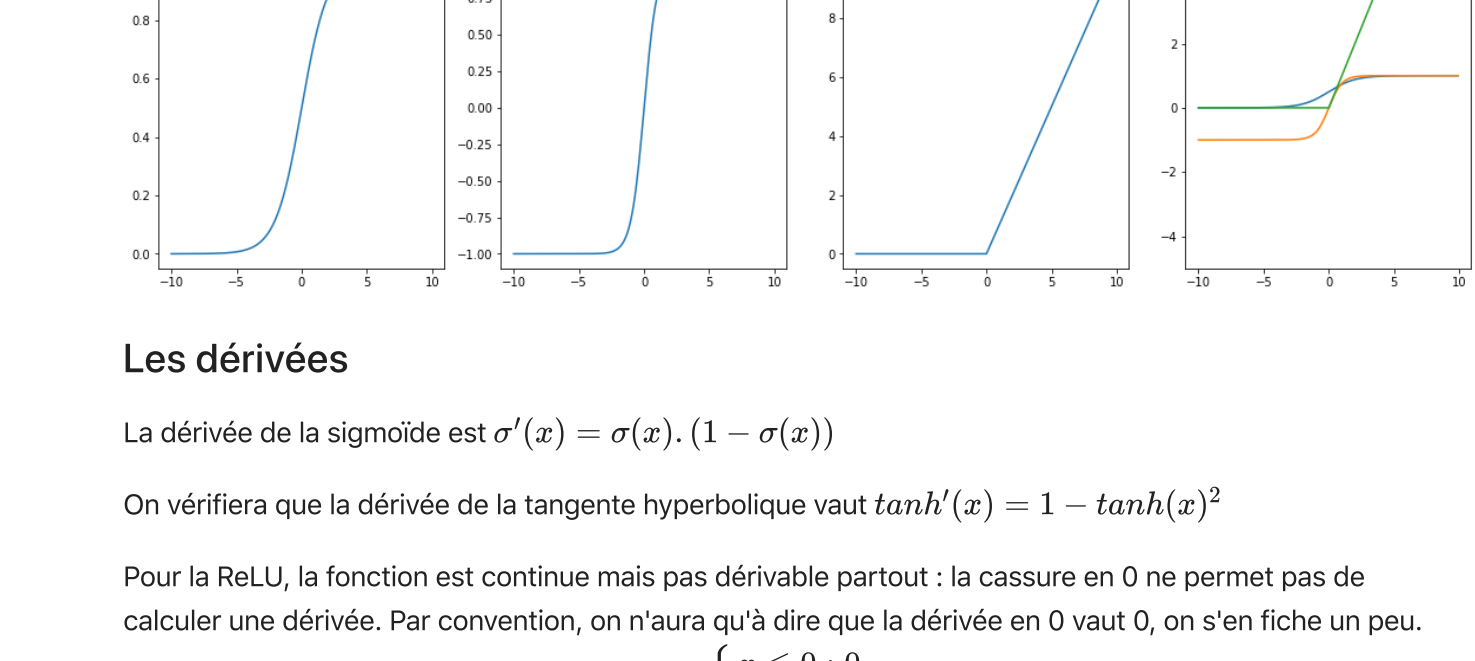
La fonction ReLU (Regularized Linear Unit) ressemble à:

- une fonction nulle pour les valeurs x<0
- une fonction identité (f(x)=x) pour x>=0

Elle est linéaire sur les deux parties, mais la cassure en 0 est suffisante pour assurer un comportement non linéaire au final.

Note: on peut utiliser au lieu de l'identité n'importe quelle fonction linéaire, mais bon en vrai on s'en fiche : si je renvoie 2x au lieu de x, les poids des neurones qui consomment mon ReLU seront divisés par 2 automatiquement quand on fera la descente de gradient alors autant rester simple.

### Les trois fonctions



### Les dérivées

La dérivée de la sigmoïde est  $\sigma'(x) = \sigma(x) \cdot (1 - \sigma(x))$

On vérifiera que la dérivée de la tangente hyperbolique vaut  $\tanh'(x) = 1 - \tanh(x)^2$

Pour la ReLU, la fonction est continue mais pas dérivable partout : la cassure en 0 ne permet pas de calculer une dérivée. Par convention, on n'aura qu'à dire que la dérivée en 0 vaut 0, on s'en fiche un peu.

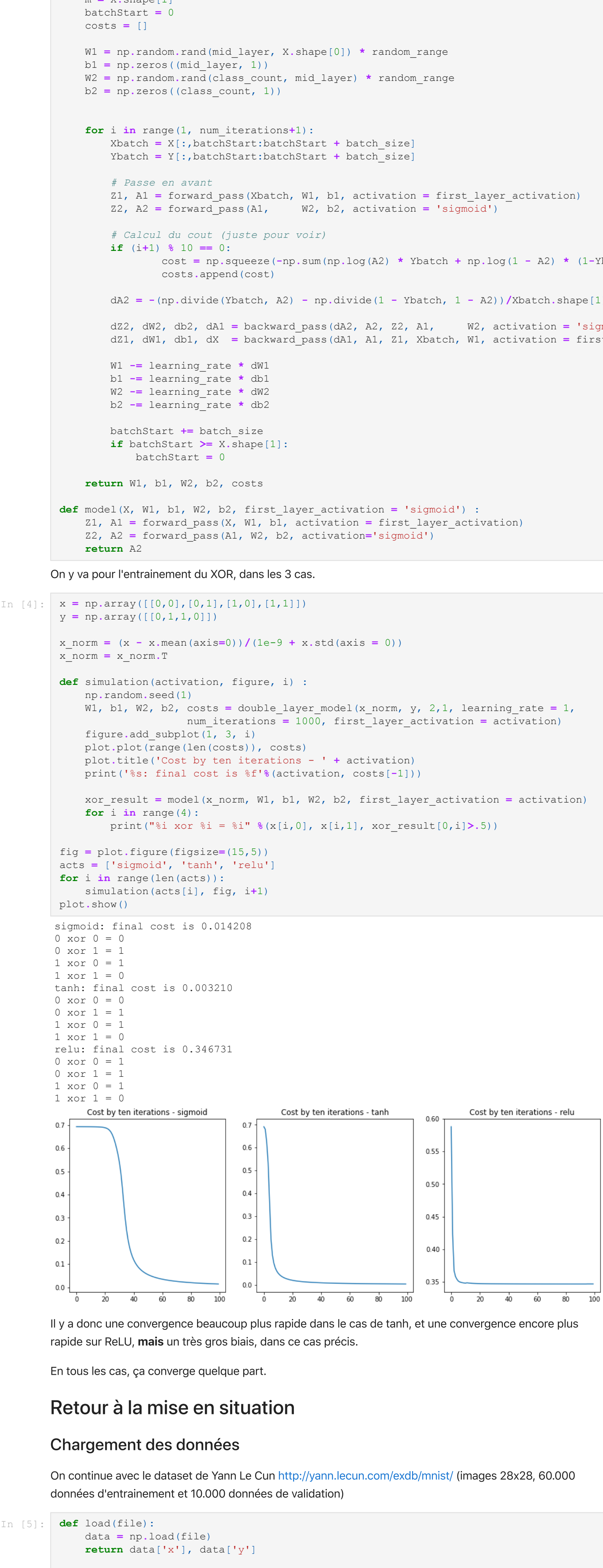
Et donc au final on aura la dérivée  $relu'(x) = \begin{cases} x \leq 0 : 0 \\ x > 0 : 1 \end{cases}$

Un gros avantage du ReLU sur les deux autres : dans le sens avant, ReLU n'est qu'un simple filtre "0 ou x", et en marche arrière sa dérivée est "0 ou 1". C'est une fonction vachement rapide à calculer, comparée aux deux autres avec des exp et des divisions dans tous les sens !

On va maintenant remplacer la première couche du réseau de neurones d'hier par la même avec une fonction d'activation différente. Pour la seconde, on garde une sigmoïde, pour tomber entre 0 et 1 c'est encore le mieux.

## Simulation rapide

Validons rapidement le concept avec une comparaison du XOR dans les 3 cas. Une précision : pour éviter des débordements, on va introduire un nouveau facteur qui indique l'amplitude du choix aléatoire initial (on génère des réels entre 0 et ce facteur qu'on). Il est crucial pour converger rapidement, et dépend beaucoup de la dimension du réseau. On verra dans la semaine comment bien choisir le facteur en question. Mais pour faire court: ReLU peut donner une valeur aussi grande que l'on veut, elle n'est pas majorée par 1, et du coup une sigmoïde appliquée par dessus donnera un résultat inintéressant proche de 0 ou de 1 (suivant le signe du W correspondant), mais enfin bref on risque de se retrouver avec une division par zéro si on dépasse la précision autorisée par Python. Et pour résoudre (en partie) le problème, on peut partir de poids beaucoup plus petits.



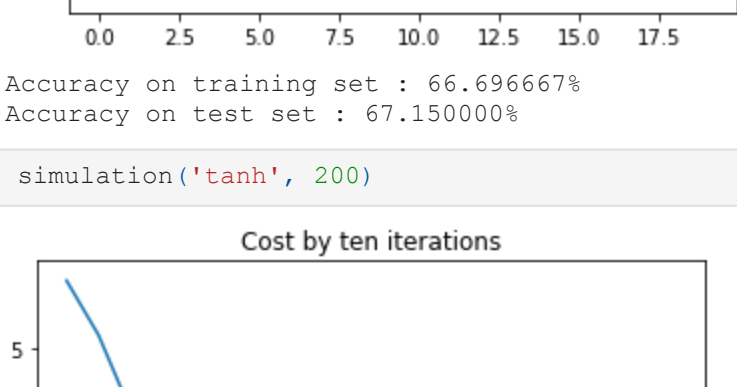
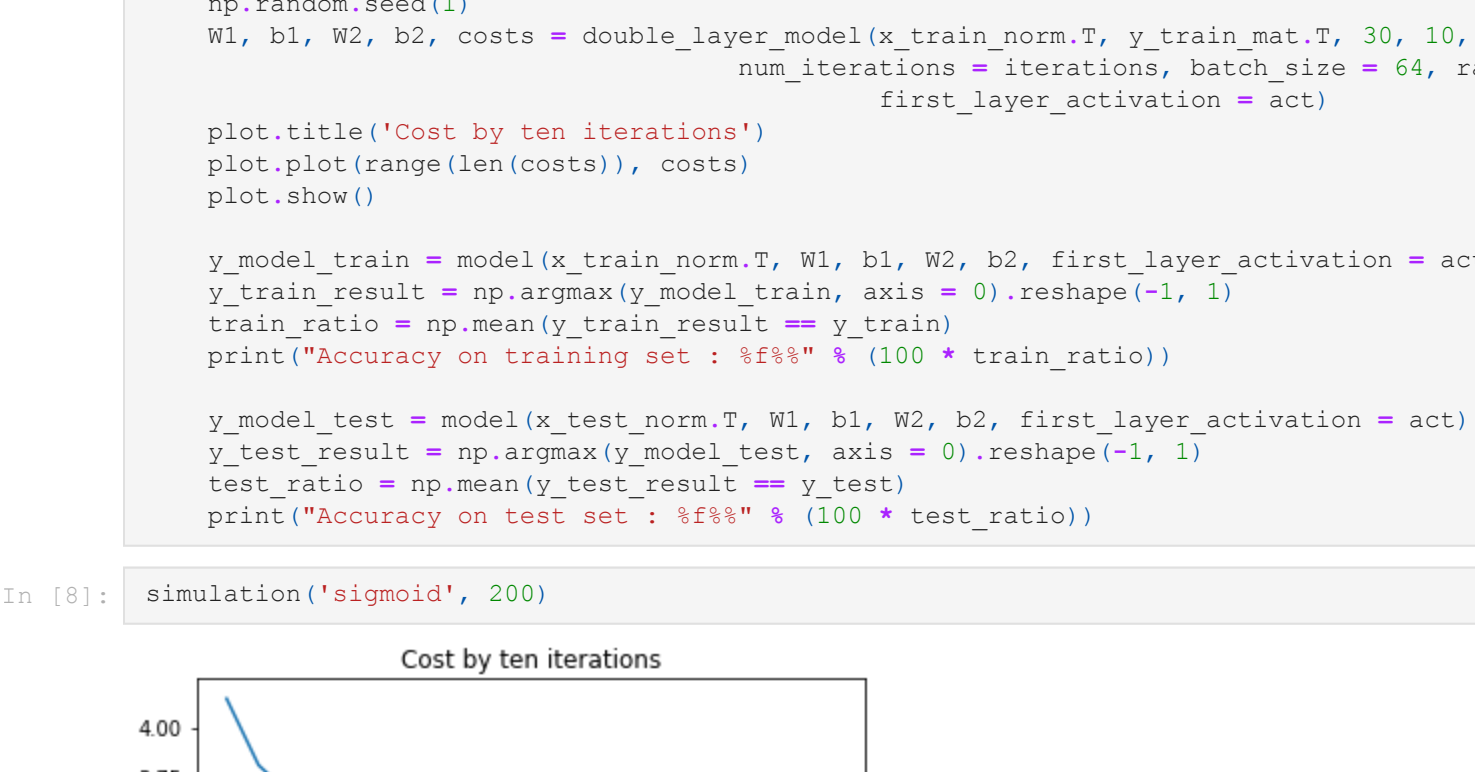
Il y a donc une convergence beaucoup plus rapide dans le cas de tanh, et une convergence encore plus rapide sur ReLU, mais un très gros biais, dans ce cas précis.

En tous les cas, ça converge quelque part.

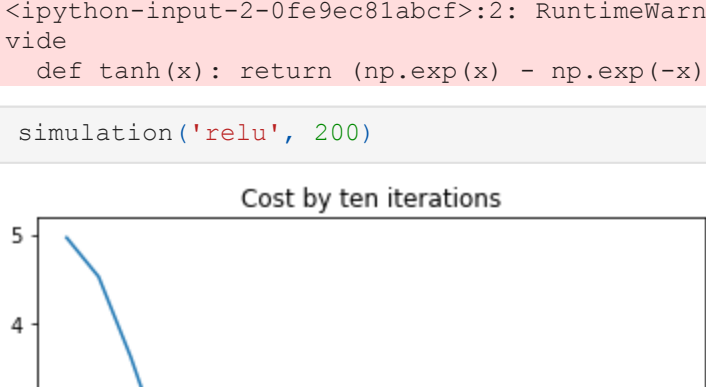
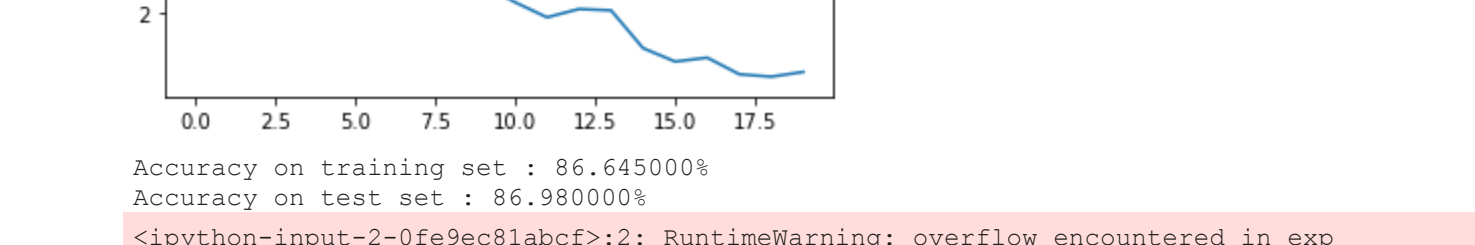
## Retour à la mise en situation

### Chargement des données

On continue avec le dataset de Yann Le Cun <http://yann.lecun.com/exdb/mnist/> (images 28x28, 60.000 données d'entrainement et 10.000 données de validation)



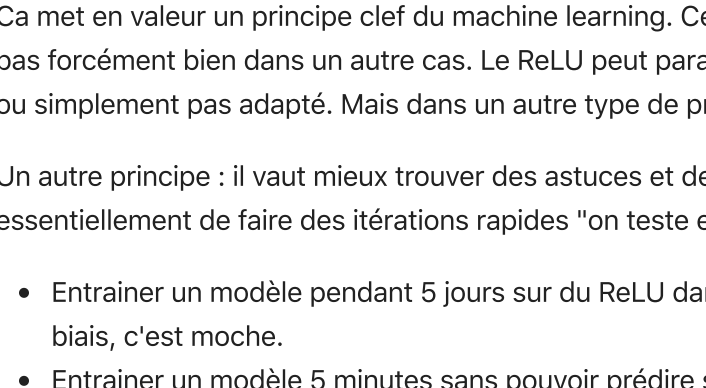
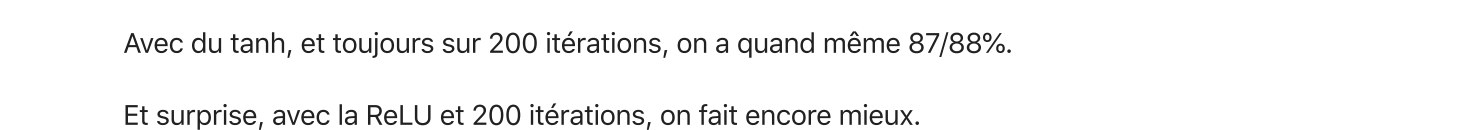
Accuracy on training set : 66.696667%  
Accuracy on test set : 67.150000%



Accuracy on training set : 86.645000%  
Accuracy on test set : 86.980000%

```
<ipython-input-2-0fe9ec81abcf>:2: RuntimeWarning: overflow encountered in exp
def tanh(x): return (np.exp(x) - np.exp(-x))/(np.exp(x) + np.exp(-x))
<ipython-input-2-0fe9ec81abcf>:2: RuntimeWarning: invalid value encountered in true_divide
```

```
def tanh(x): return (np.exp(x) - np.exp(-x))/(np.exp(x) + np.exp(-x))
```



Accuracy on training set : 89.370000%  
Accuracy on test set : 89.670000%

```
<ipython-input-2-0fe9ec81abcf>:1: RuntimeWarning: overflow encountered in exp
def sigmoid(x): return 1 / (1 + np.exp(-x))
```

On remarque que la sigmoïde, qui donnait de bons résultats hier, est en fait pas si géniale que ça : ok, ça converge. Mais avec 200 itérations, on n'a pas un résultat exceptionnel...

Avec du tanh, et toujours sur 200 itérations, on a quand même 87/88%.

Et surprise, avec la ReLU et 200 itérations, on fait encore mieux.

Ca met en valeur un principe clef du machine learning. Ce qui marche bien sur un modèle ne marchera pas forcément bien dans un autre cas. Le ReLU peut paraître mauvais comparé à un traitement du XOR, ou simplement pas adapté. Mais dans un autre type de problème il sera efficace.

Un autre principe : il vaut mieux trouver des astuces et des conditions de convergence **rapide**. Ca permet essentiellement de faire des itérations rapides "on teste et on avise". Exemples et contre-exemples:

- Entraîner un modèle pendant 5 jours sur du ReLU dans le cas du XOR pour voir qu'on a un souci de biais, c'est moche.
- Entraîner un modèle 5 minutes sans pouvoir prédire son comportement, c'est inutile.
- Entraîner de la reconnaissance de chiffres sur peu d'itérations et plusieurs fonctions d'activations, pour valider que ReLU se comporte mieux, permet ensuite lancer un entraînement beaucoup plus long.

Le modèle n'est pas forcément meilleur, mais on converge plus vite. Et comme on l'a dit, un gros avantage de ReLU est d'être calculable très rapidement aussi, donc chaque itération prendra en plus moins de temps! Dans des problèmes beaucoup plus complexes, la convergence ne s'atteint pas en quelques secondes : du coup toute technique qui permet de converger plus vite est bonne à prendre ! Et on en verra encore d'autres.