# Algorithmes, Programmation, IA

Cours 10

Jean-Jacques Lévy

jean-jacques.levy@inria.fr

### Plan

- apprentissage profond
- caculs des poids (simple descente du gradient, vectorisation, Tensorflow)
- exemple de traitement de langue naturelle
- divers types de modèles

# Python++ (fonctions)

• les arguments des fonctions peuvent être étiquetés

```
def address (nom, rue, tel):
    print (nom + " " + rue); print (tel)

address (rue = "Couedic", nom = "JJL", tel = 3888)
```

• les arguments étiquetés peuvent avoir une valeur par défaut

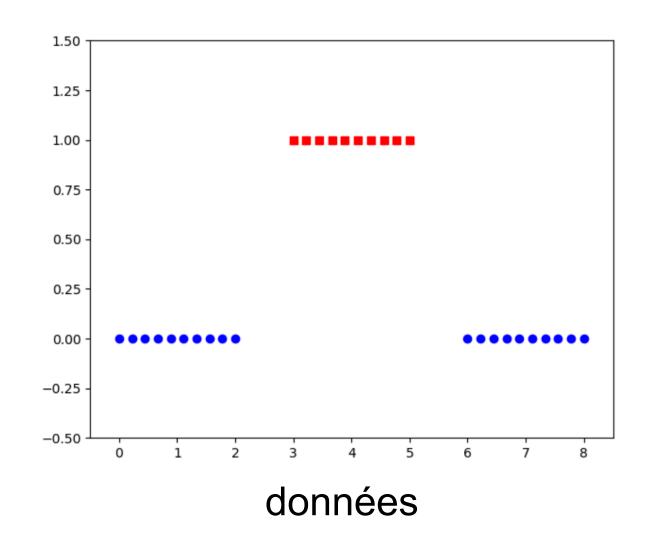
```
def address (nom, rue, tel='3523') :
    print (nom + " " + rue); print (tel)

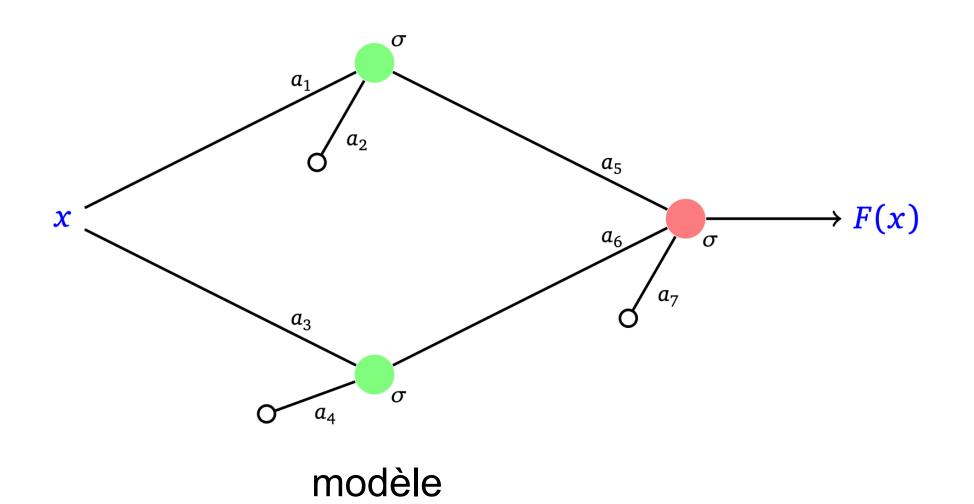
address (rue = "Couedic", nom = "JJL")
```

• les arguments étiquetés doivent apparaître après les arguments positionnels

# Calcul des poids

• on dispose d'un jeu de données d'entrainement (training set)





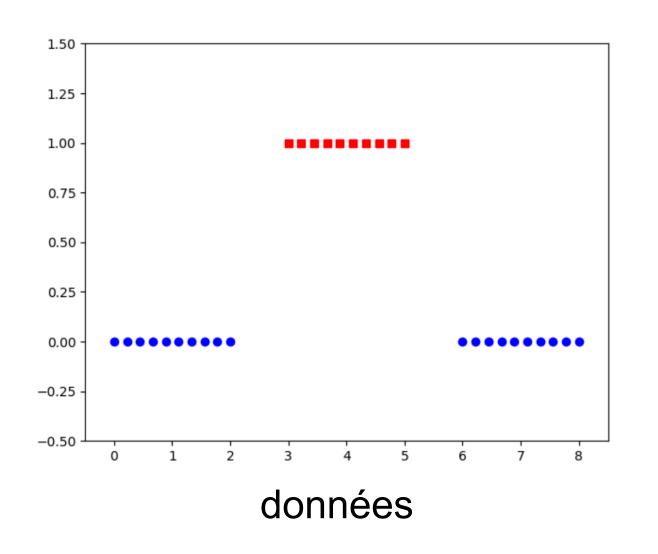
Essayer une descente de gradient directement sur les 7 poids du réseau.

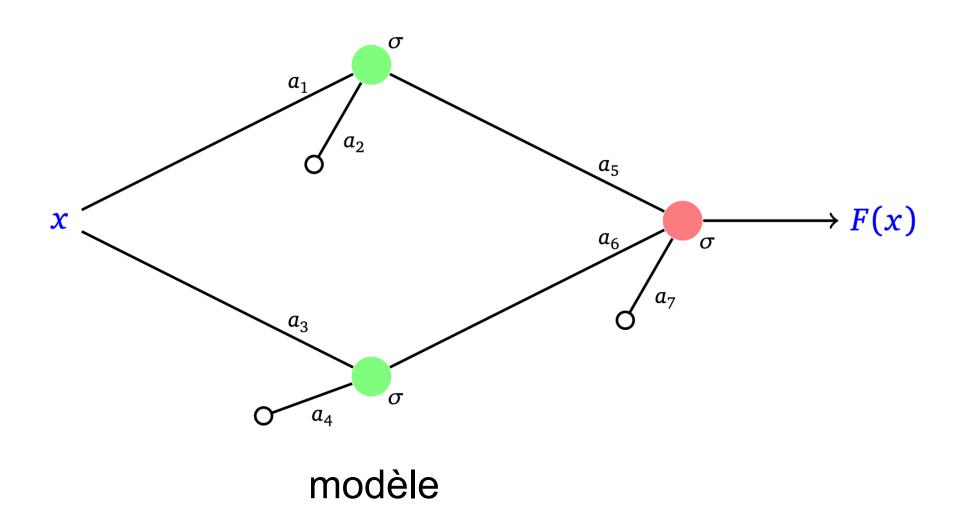
Essayer ensuite avec rétro-propagation..

Essayer encore plus tard avec Pytorch.

# Rappel

• on dispose d'un jeu de données d'entrainement (training set)





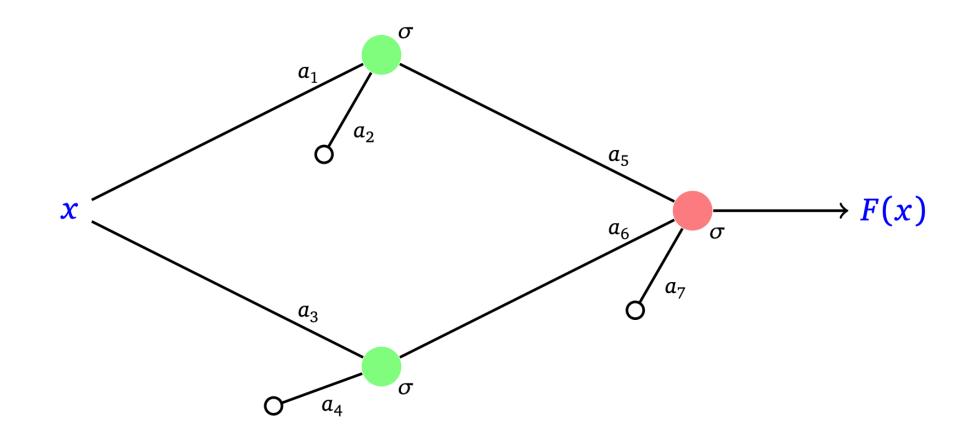
- on veut trouver un modèle (neural net) qui approxime les données d'entrainement
- ce qui permet de prevoir des résultats en dehors des données d'entrainement (predictions)
- dans notre exercice, on suppose un réseau à 3 couches (entrée, intermédiaire avec 2 neurones, sortie) et une fonction sigmoïde d'activation
- il reste à calculer les paramètres (poids)

• jeu de données et coefficients initiaux

```
bleus = [(i/10, 0.2) for i in range (20)] + \
        [(6 + i/10, 0.2) for i in range (20)]
rouges = [(3 + i/10, 1) for i in range (20)]

dataset = bleus + rouges

# --- paramètres (poids) du réseau ---
a1, a2, a3, a4 = (0.0, 1.0, 0.0, -1.0)
a5, a6, a7 = (1.0, 1.0, -1.0)
```



fonction d'activation (et sa dérivée)

```
def sigma (x):
    return 1 / (1 + np.exp(-x))

def dsigma (x):
    return (sigma (x)) * (1 - sigma(x))
```

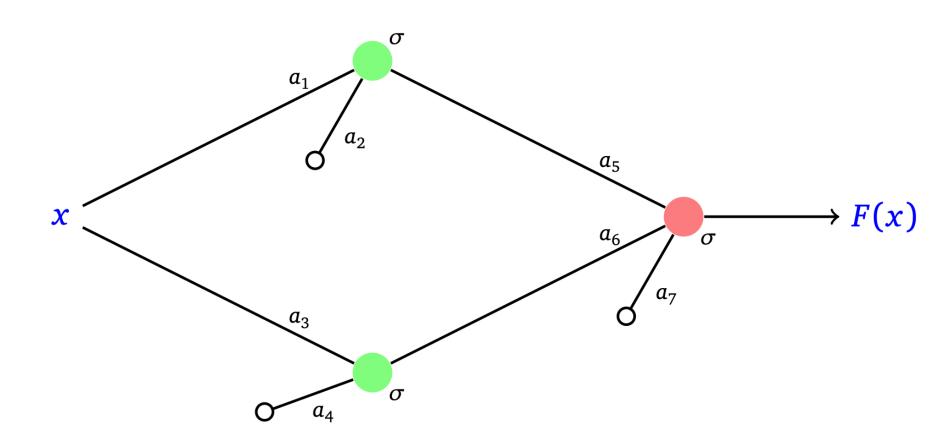
• on va calculer les poids par une descente de gradient sur les 7 coefficients

fonction calculée

```
def F1(x):
    global a1, a2, a3, a4, a5, a6, a7
    return sigma (a1 * x + a2)

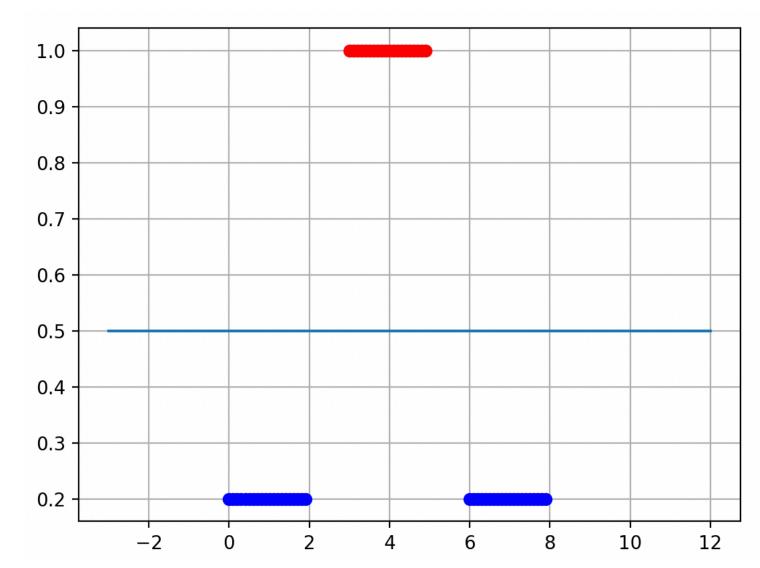
def F2(x):
    global a1, a2, a3, a4, a5, a6, a7
    return sigma (a3 * x + a4)

def F(x):
    global a1, a2, a3, a4, a5, a6, a7
    return sigma (a5 * F1(x) + a6 * F2(x) + a7)
```



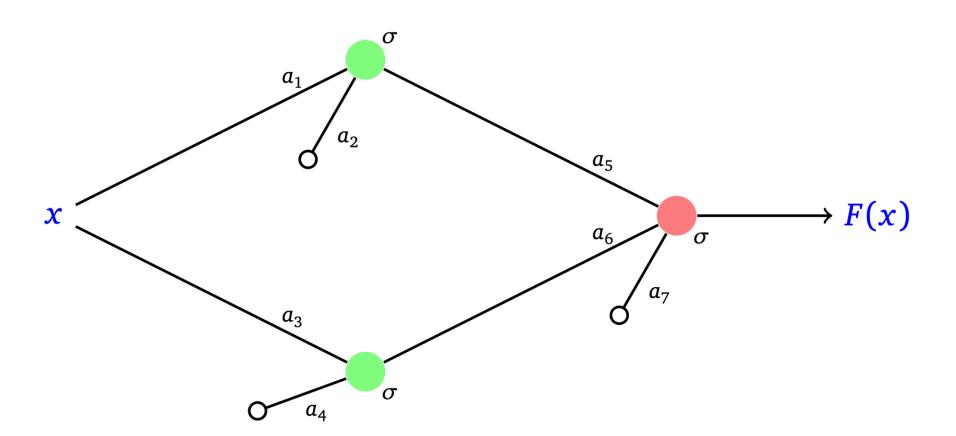
affichage des données et de la fonction calculée

```
xlim = (-3, 12)
def show (dataset):
    global xlim
    for (x, y) in dataset:
        plt.plot (x, y, marker='o', color='r' if y == 1 else 'b')
    X = np.linspace (*xlim, 100)
    Y = F(X)
    plt.plot (X, Y)
    plt.grid()
    plt.show()
```



dérivées de la fonction calculée

```
def dF1(x):
   global a1, a2, a3, a4, a5, a6, a7
   return np.array ([dsigma (a1 * x + a2) * x, \
            dsigma (a1 \times x + a2), \
            0, 0, 0, 0, 0])
def dF2(x):
    global a1, a2, a3, a4, a5, a6, a7
    return np.array ([0, 0, dsigma (a3 * x + a4) * x,
            dsigma (a3 * x + a4), 0, 0, 0]
def dF(x):
   global a1, a2, a3, a4, a5, a6, a7
   dA3 = dsigma(a5 * F1(x) + a6 * F2(x) + a7)
   return np.array ([ \
            dA3 * a5 * dF1(x)[0],
            dA3 * a5 * dF1(x)[1],
            dA3 * a6 * dF2(x)[2],
            dA3 * a6 * dF2(x)[3],
            dA3 * \mathbf{F1}(x),
            dA3 * F2(x),
            dA3 ])
```



la fonction J de coût

```
def Jxy(x, y):
    return 0.5 * (F (x) - y)**2

def J(dataset):
    s = 0; n = len(dataset)
    for (x, y) in dataset:
        s += Jxy (x, y)
    return s
```

### dérivées du coût

```
def derJxy (x, y):
    return (F(x) - y) * dF (x)

def derJ (dataset):
    s = 0
    for (x, y) in dataset:
        s += derJxy (x, y)
    return s
```

### programme principal

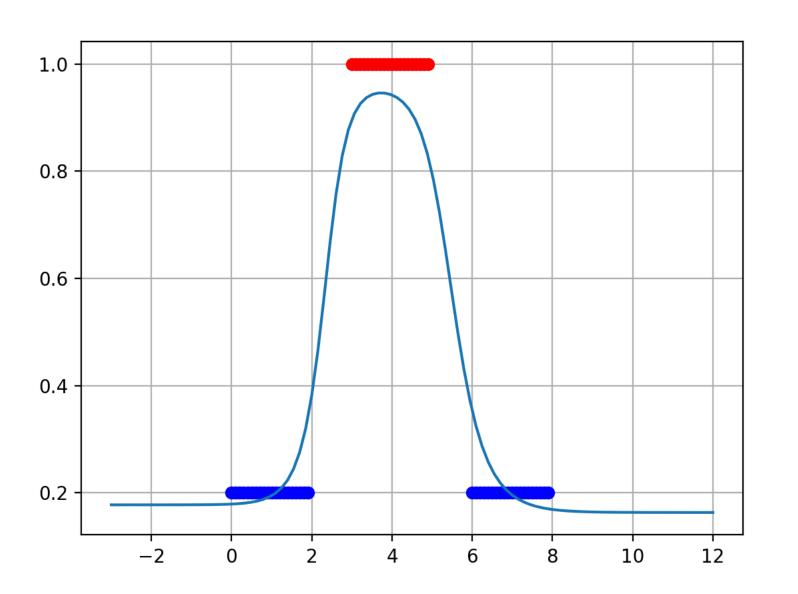
```
alpha = 0.1
print ('J = ', J (dataset))
print ('a1a2a3a4a5a6a7 --> ', a1, a2, a3, a4, a5, a6, a7)
show (dataset)

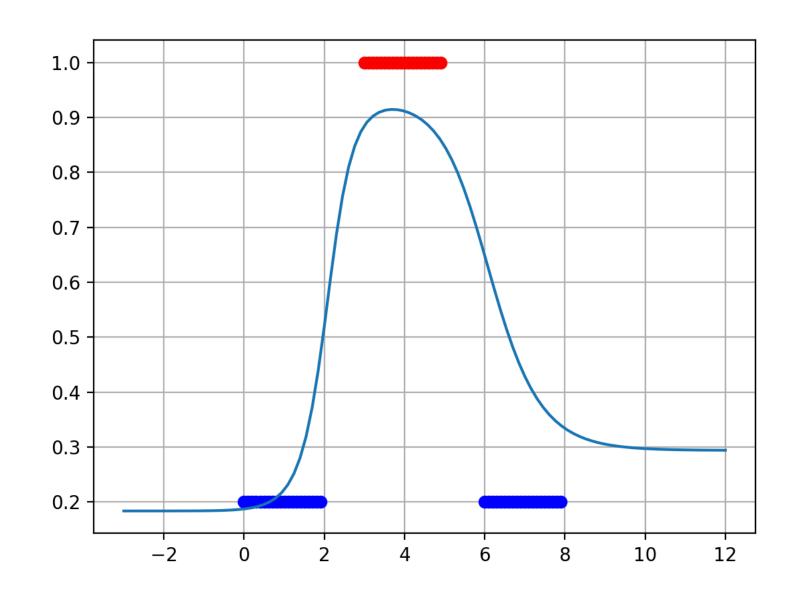
descent (dataset)

print ('J = ', J (dataset))
print ('a1a2a3a4a5a6a7 --> ', a1, a2, a3, a4, a5, a6, a7)
show (dataset)
```

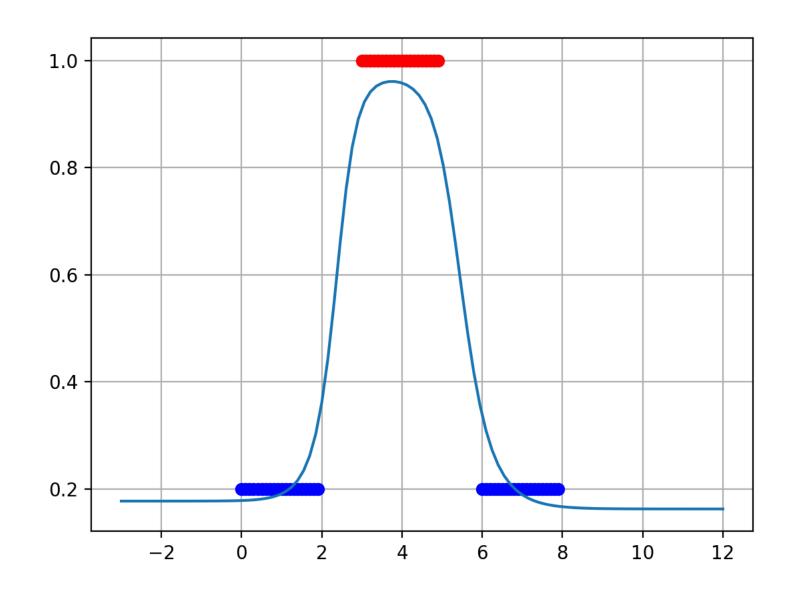
### descemte du gradient

```
def descent (dataset) :
    global a1, a2, a3, a4, a5, a6, a7
    global alpha
    for epoch in range (800) :
        dd = derJ (dataset)
        a1 -= alpha * dd [0]
        a2 -= alpha * dd [1]
        a3 -= alpha * dd [2]
        a4 -= alpha * dd [3]
        a5 -= alpha * dd [4]
        a6 -= alpha * dd [5]
        a7 -= alpha * dd [6]
```

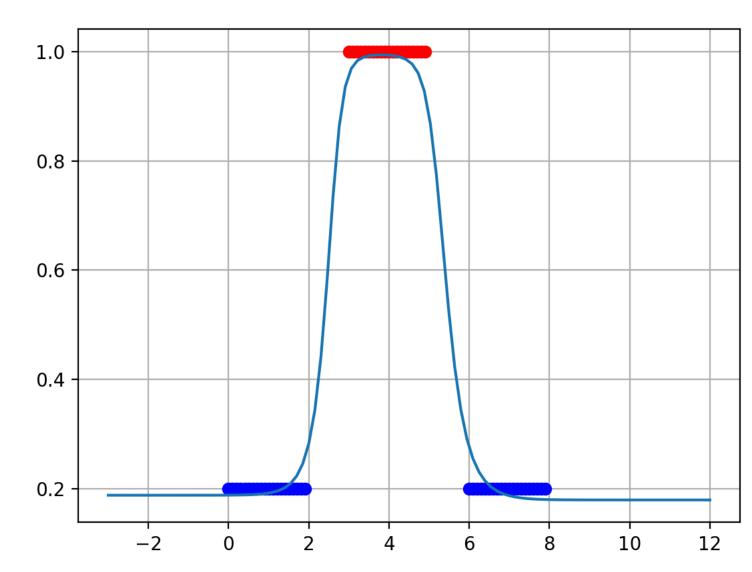








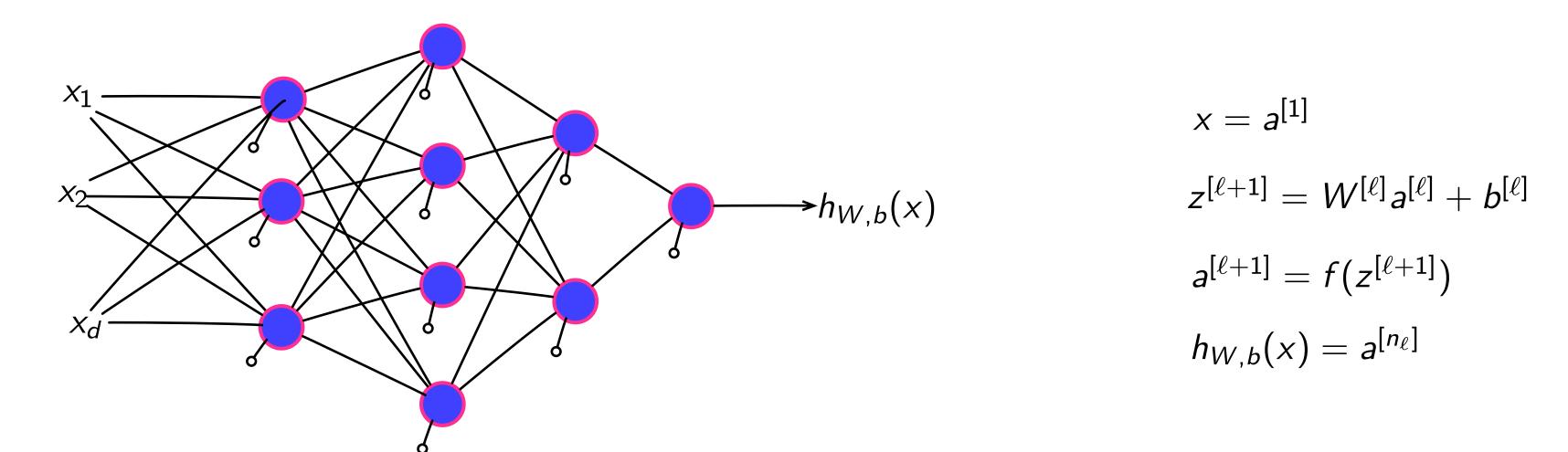




5000 itérations

# Réseau de neurones (rappel)

• exemple de réseau avec d entrées, 1 sortie, 5 couches  $(n_{\ell} = 5)$  et une fonction d'activation uniforme f



$$d = d_1$$
  $n_{\ell} = 5$   $W^{[1]}: 3 \times d$   $W^{[2]}: 4 \times 3$   $W^{[3]}: 2 \times 4$   $W^{[4]}: 1 \times 2$ 

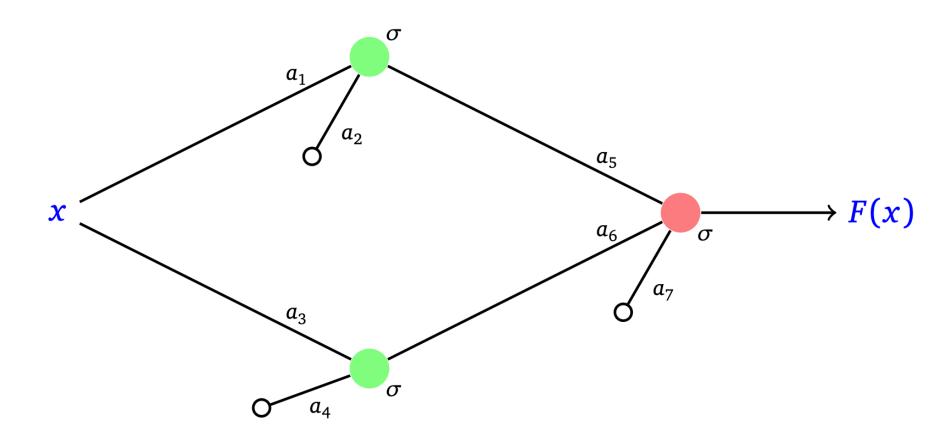
• 31 paramètres quand d = 3 !!

### [version algébrique]

fonction calculée

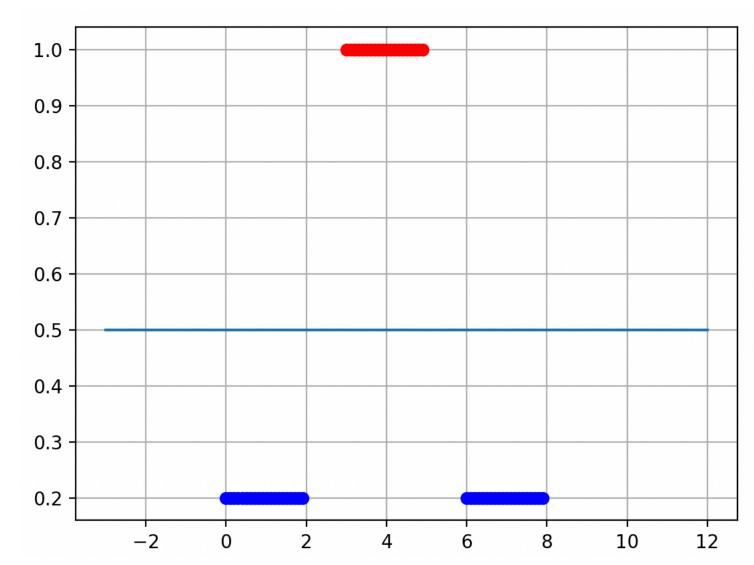
```
W1 = np.array ([[a1], [a3]]) # shape 2x1
B1 = np.array ([[a2], [a4]]) # shape 2x1
W2 = np.array ([[a5, a6]]) # shape 1x2
B2 = np.array ([[a7]]) # shape 1x1

def F(x):
    global W1, W2, B1, B2
    X = np.array([[x]])
    Z2 = np.dot (W1, X)+ B1; A2 = sigma (Z2)
    Z3 = np.dot (W2, A2) + B2; A3 = sigma (Z3)
    return A3[0,0]
```



• affichage des données et de la fonction calculée

```
xlim = (-3, 12)
def show (dataset):
    global xlim
    for (x, y) in dataset:
        plt.plot (x, y, marker='o', color='r' if y == 1 else 'b')
    X = np.linspace (*xlim, 100)
    Y = np.array ([F(x) for x in X])
    plt.plot (X, Y)
    plt.grid()
    plt.show()
```



# Rétro-propagation (rappel)

- calcul des dérivées partielles par rapport aux poids
  - 1) on calcule les  $a_i^{[\ell]}$  et les  $z_i^{[\ell]}$  par une passe en avant
  - 2) pour la sortie (ou les sorties si plusieurs), on calcule

$$\delta^{[n_\ell]} = (a^{[n_\ell]} - y) \bullet f'(z^{[n_\ell]})$$

3) pour chaque neurone des couches intermédiaires, on calcule

$$\delta^{[\ell]} = ((\mathcal{W}^{[\ell]})^{\mathsf{T}} \ \delta^{[\ell+1]}) \bullet f'(z^{[\ell]})$$

4) les dérivées sont maintenant établies

$$\nabla_{W^{[\ell]}}J(W,b;x,y)=\delta^{[\ell+1]}(a^{[\ell]})^T$$

$$\nabla_{b^{[\ell]}}J(W,b;x,y)=\delta^{[\ell+1]}$$

multiplication point par point

 ○ opérateur de dérivation par rapport à tous les éléments de la matrice

• fonction de coût

```
def Jxy(x, y):
    return 0.5 * (F (x) - y)**2

def J(dataset):
    s = 0
    for (x, y) in dataset:
        s += Jxy (x, y)
    return s
```

programme principal

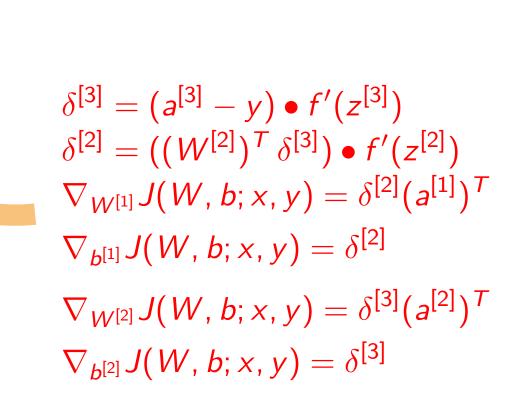
```
alpha = 0.1
print ('J = ', J (dataset))
print ('w1b2w2b2 ---> ', W1, B1, W2, B2)
show(dataset)
descent (dataset)
print ('J = ', J (dataset))
print ('w1b2w2b2 ---> ', W1, B1, W2, B2)
show(dataset)
```

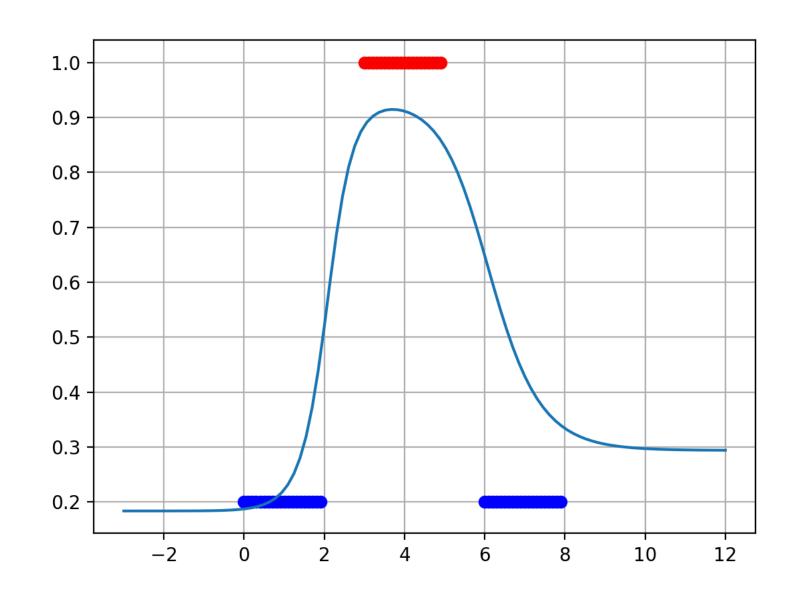
• dérivée du coût

```
def derJxy (x, y):
   # forward pass
   global W1, W2, B1, B2
   X = np.array([[x]])
   Z2 = np.dot (W1, X) + B1; A2 = sigma (Z2)
   Z3 = np.dot (W2, A2) + B2; A3 = sigma (Z3)
   # back propagate
   Y = np.array([[y]])
   D3 = (A3 - Y) * dsigma(Z3)
   D2 = np.dot (W2.T, D3) * dsigma (Z2)
   dJW1 = np.dot (D2, X.T); dJB1 = D2
   dJW2 = np.dot (D3, A2.T); dJB2 = D3
   return dJW1, dJW2, dJB1, dJB2
def derJ (dataset) :
   sW1, sW2, sb1, sb2 = 0, 0, 0
   for (x, y) in dataset :
        (dW1, dW2, dB1, dB2) = derJxy(x, y)
       sW1 += dW1; sW2 += dW2
       sb1 += dB1; sb2 += dB2
   return (sW1, sW2, sb1, sb2)
```

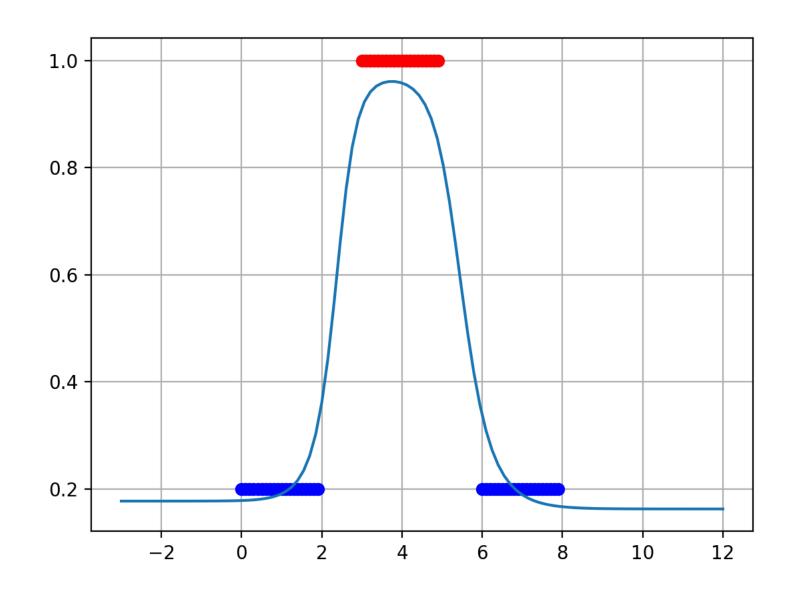
### descente du gradient

```
def descent (dataset) :
    global W1, W2, B1, B2
    global alpha
    for epoch in range (1000) :
        (dW1, dW2, dB1, dB2) = derJ (dataset)
        W1 -= alpha * dW1; B1 -= alpha * dB1
        W2 -= alpha * dW2; B2 -= alpha * dB2
```

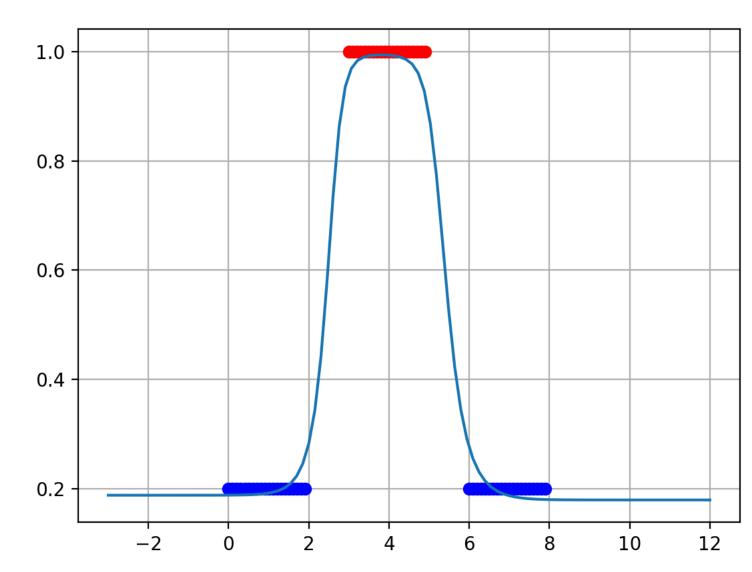












5000 itérations

### Apprentissage profond - avec Keras

• importation des modules au dessus de Tensorflow [ google ] et Keras [François Chollet, 2015]

```
from tensorflow import keras
from keras models import Sequential
from keras layers import Dense
from keras import optimizers
```

• fonction à approcher

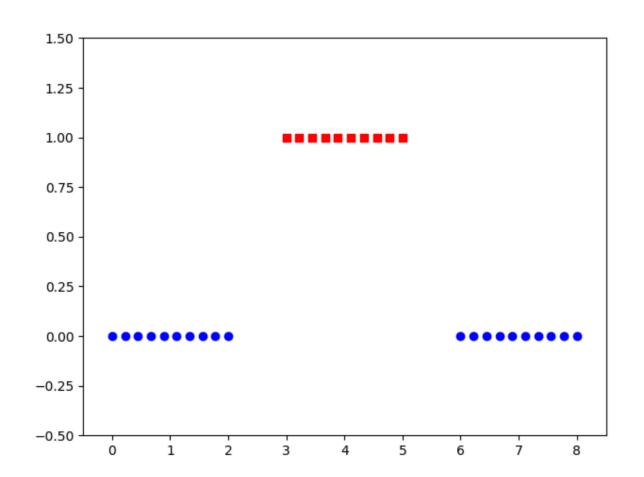
```
def F(x):
return 1 if 3 <= x and x <= 5 else 0
```

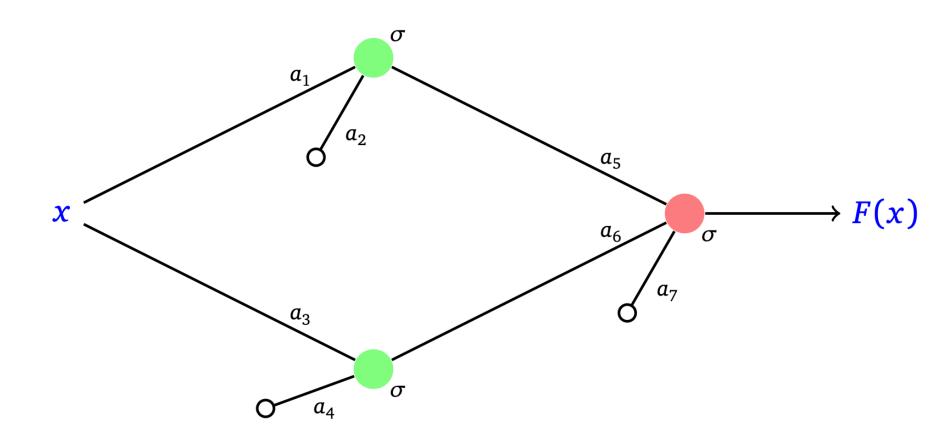
• jeu de données et coefficients initiaux

```
a, b = -2, 9

N = 100
X = np.linspace(a, b, N)
Y = np.array ([F(x) for x in X])

X_train = X.reshape(-1,1)
Y_train = Y.reshape(-1,1)
```





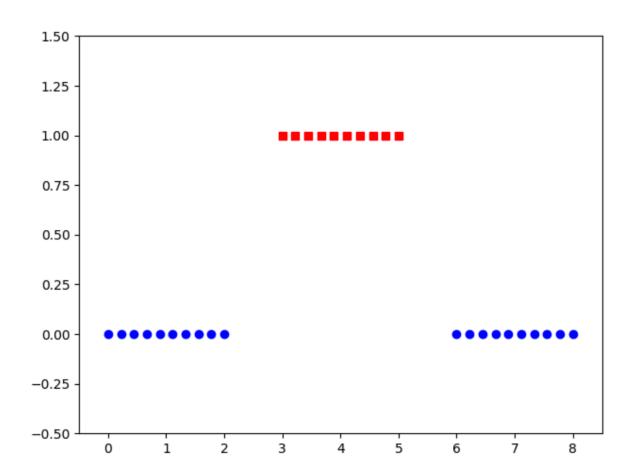
### Apprentissage profond - avec Keras

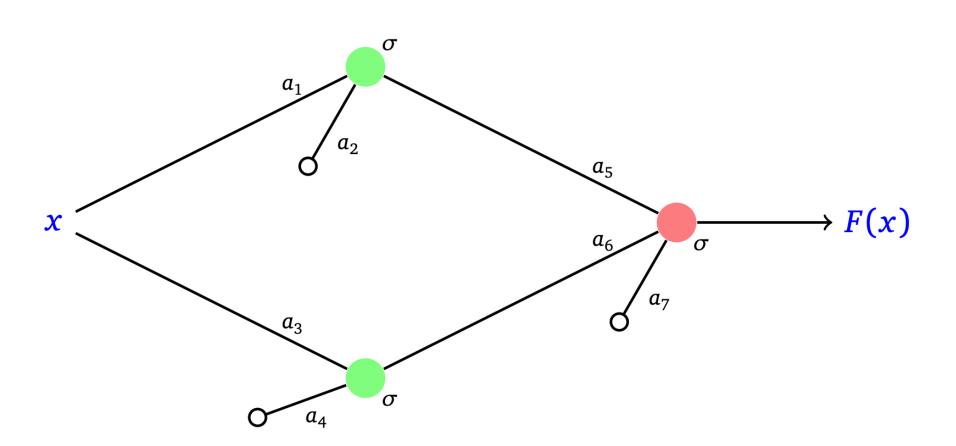
### définition du réseau

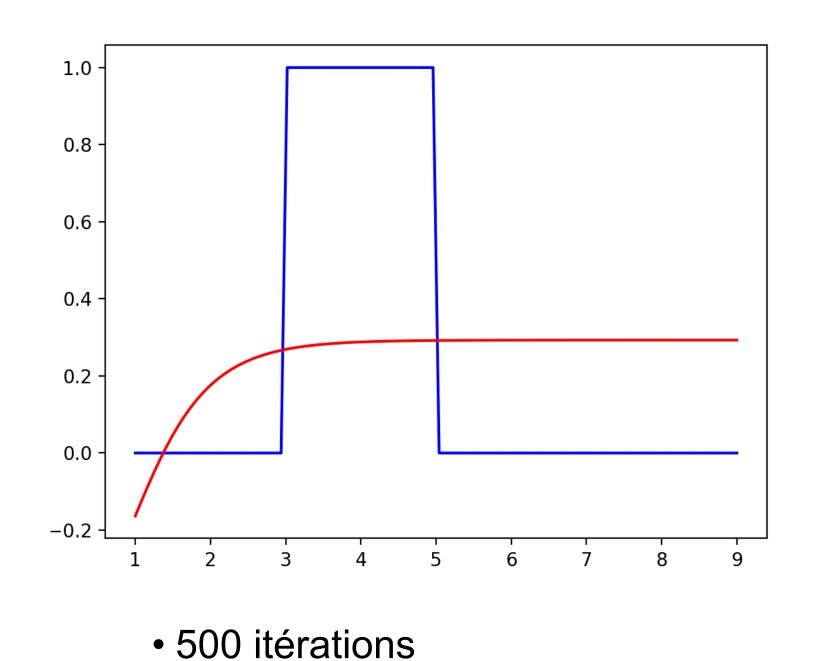
### apprentissage

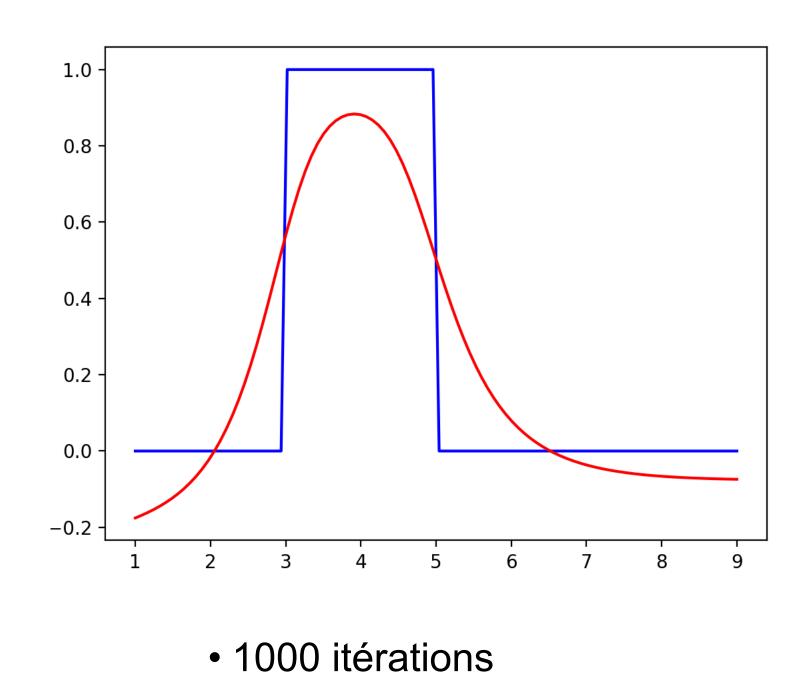
### visualisation du résultat

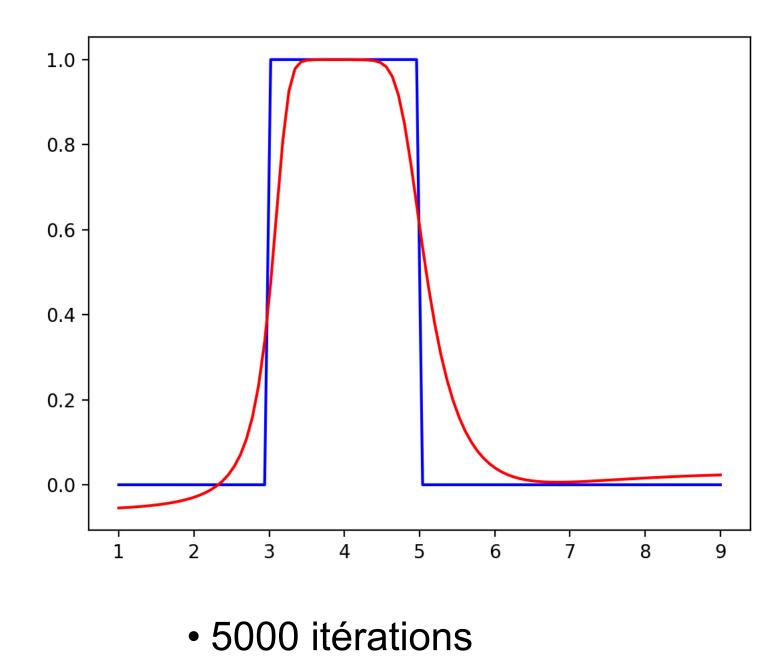
```
Y_predict = modele.predict(X_train)
plt.plot(X_train, Y_train, color='blue')
plt.plot(X_train, Y_predict, color='red')
plt.show()
```











Exercice Modifier le modèle pour de meilleurs résultats.

Exercice Essayer avec d'autres fonctions plus dérivables.

- l'initialisation des poids est importante pour diminuer l'erreur
- beaucoup d'opérations peuvent être vectorisées avec un GPU (processeur graphique)
  - la multiplication de matrices exécutée en parallèle [https://youtu.be/cGEIEnekmRM]
  - l'addition de matrices exécutée en parallèle
  - chaque étape de la rétro-propagation en parallèle
  - tout calcul sans variables modifiables (ce n'est pas le cas de la descente de gradient)
- NVDIA est actuellement le plus important fabricant de GPU, aussi rebaptisé Neural Engine chez Apple
- Tensorflow, Pytorch marchent avec GPU sur PC avec chip AMD, macbook-pro, et iMac avec les chips M1, .. , M4 (chip Apple/ARM)

### Reconnaissance de texte

[ exemple pris de Deepmath Exo7 (cité dans la 2ème diapositive), page 192 ]

Le but de cet exemple est de décider si une critique de film est positive ou négative. Voici un exemple de critique (la critique numéro 123) de la base que nous utiliserons :

beautiful and touching movie rich colors great settings good acting and one of the most charming movies i have seen in a while i never saw such an interesting setting when i was in china my wife liked it so much she asked me to log on and rate it so other would enjoy too

C'est un critique positive!

La base IMDB fournit 25000 critiques d'apprentissage, chacune étant déjà catégorisée positive (valeur 1) ou négative (valeur 0)

Voici trois (fausses) critiques :

- « bon film à voir absolument»
- « début moyen mais après c'est très mauvais»
- « film drôle et émouvant on passe un bon moment»

• Tout d'abord, parmi toutes les critiques, on ne retient que les 10 mots les plus fréquents. Imaginons que ce sont les mots : film, bon, mauvais, voir, éviter, très, bien, mal, moyen, super

En réalité on retiendra les 1000 (voire 10 000) mots les plus fréquents.

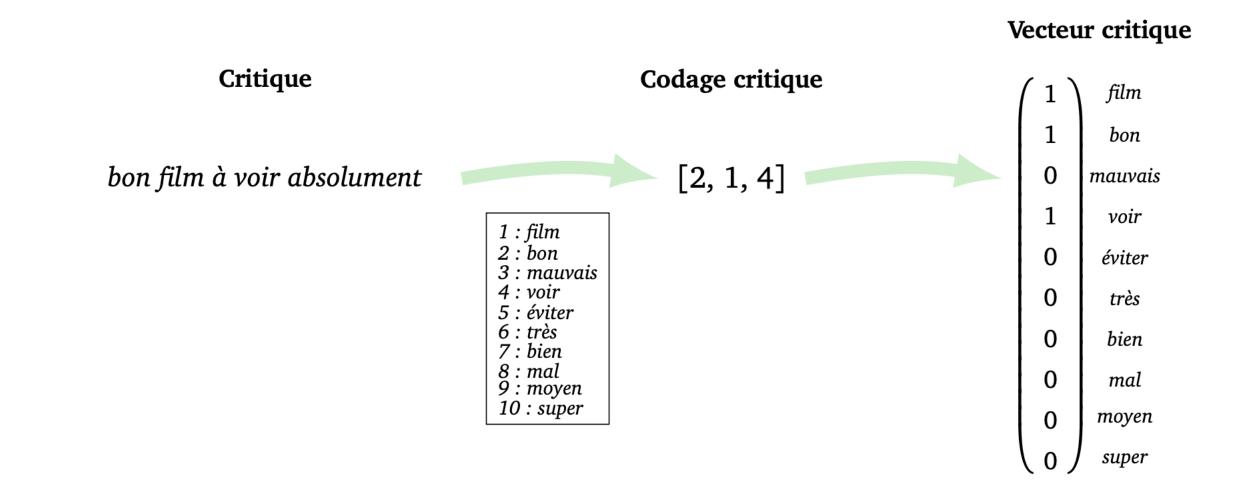
### Reconnaissance de texte

[ exemple pris de Deepmath (cité dans la 2ème diapositive), page 192 ]

• Tout d'abord, parmi toutes les critiques, on ne retient que les 10 mots les plus fréquents. Imaginons que ce sont les mots : film, bon, mauvais, voir, éviter, très, bien, mal, moyen, super

En réalité on retiendra les 1000 (voire 10 000) mots les plus fréquents.

- On code une phrase comme une liste d'indices de mots. Le mot « film » est remplacé par 1, le mot « bon » par 2,... jusqu'à « super ». Les autres mots ne sont pas pris en compte. Par exemple la critique « bon film à voir absolument » devient la liste d'indices [2, 1, 4]. Les deux autres critiques deviennent [9, 3] et [1, 2].
- On transforme ensuite chaque liste en un vecteur de taille fixe n = 10 de coordonnées 0 ou 1. On place un 1 en position i si le mot numéro i apparaît dans la critique. Ainsi la phrase « bon film à voir absolument », codée en [2, 1, 4] devient le vecteur X = (1, 1, 0, 1, 0, 0, 0, 0, 0, 0). La phrase codée [9, 3] donne le vecteur X'= (0, 0, 1, 0, 0, 0, 0, 0, 1, 0).



# Apprentissage supervisé, réinforcé ou non

- apprentissage supervisé avec un jeu de données d'entrainement où la sortie est étiquetée
  - ce qu'on a vu jusqu'à présent

- apprentissage non supervisé avec un jeu de données non étiquetées
  - on essaie de les classer en groupes à partir de données prises au hasard initialement

- apprentissage renforcé est un mélange des deux techniques
  - on modifie l'apprentissage automatique avec quelques interventions humaines
- l'apprentissage génératif désigne la prédiction à partir d'une des méthodes précédentes
  - par exemple, GPT (generative pre-training transformer) utilisé dans les LLM ou Chat-GPT

### Modèles de réseaux

- nn: modèles séquentiels
  - suite de couches (plus ou moins denses, couches de convolution possibles)

- rnn: modèles récurrents
  - boucles possibles dans l'enchaînement des couches

- modèles fonctionnels (en Keras)
  - topologie à la carte

### Conclusion

- l'apprentissage est bien loin des algorithmes
- la contribution informatique est : le stockage possible de grand nombre de données, les processeurs parallèles
- l'apprentissage est une discipline plutôt expérimentale (statistiques et probabilités)
- l'apprentissage a été souvent le domaine de la vision par ordinateur (computer vision)
- de manière surprenante, l'apprentissage marche pour l'analyse de la langue naturelle (NLP) avec les LLM et GPT

- il reste à l'intégrer à la démonstration automatique de théorèmes et à la vérification du logiciel
- etc..
- prochain cours: algorithmes géométriques