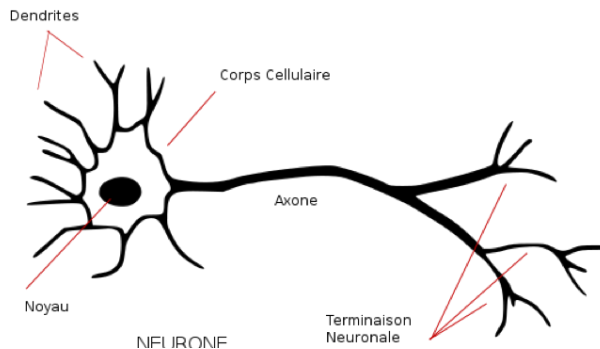


## Les réseaux de neurones

### Les neurones

"Un neurone, ou une cellule nerveuse, est une cellule excitable constituant l'unité fonctionnelle de base du système nerveux. Les neurones assurent la transmission d'un signal bioélectrique appelé influx nerveux. Ils ont deux propriétés physiologiques : l'excitabilité, c'est-à-dire la capacité de répondre aux stimulations et de convertir celles-ci en impulsions nerveuses, et la conductivité, c'est-à-dire la capacité de transmettre les impulsions." (Wikipedia)

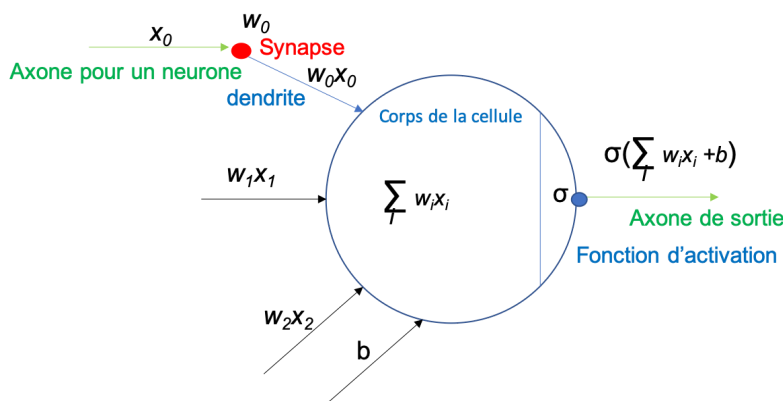
La structure d'un neurone (source : [https://fr.wikipedia.org/wiki/Fichier:Neurone\\_-\\_commenté.svg](https://fr.wikipedia.org/wiki/Fichier:Neurone_-_commenté.svg) ([https://fr.wikipedia.org/wiki/Fichier:Neurone\\_-\\_commenté.svg](https://fr.wikipedia.org/wiki/Fichier:Neurone_-_commenté.svg))) :



Le fonctionnement est le suivant : tout d'abord, les dendrites reçoivent l'influx nerveux d'autres neurones. Le neurone évalue alors l'ensemble de la stimulation reçue. Si celle-ci est suffisante, il est excité : il transmet un signal (0/1) le long de l'axone et l'excitation est propagée jusqu'aux autres neurones qui y sont connectés via les synapses.

"Un réseau de neurones artificiels, ou réseau neuronal artificiel, est un système dont la conception est à l'origine schématiquement inspirée du fonctionnement des neurones biologiques" (Wikipedia)

La structure d'un neurone artificiel :



Comme nous le constatons, un neurone artificiel est assez similaire à un neurone. Il comprend un ensemble d'entrées (synapses) auxquelles un ensemble de poids sont ajoutés (dans le notebook sur la descente de gradient, ces poids correspondent aux paramètres qu'il fallait trouver pour les fonctions linéaires -  $\theta$ ). Il possède également une entrée particulière appelée *biais*. Une fonction additive (combinaison linéaire) calcule la somme pondérée des entrées :  $\sum_i w_i x_i$ . La sortie du noeud est déterminée en appliquant une fonction de transfert non-linéaire,  $\sigma(\sum_i w_i x_i + b)$ .

## Lorsque la régression logistique ne fonctionne plus

Dans le notebook sur la descente de gradient, nous avons terminé par la régression logistique et avons vu qu'il était possible d'afficher les limites de décision (une droite) pour classer les iris. Considérons, à présent, la figure suivante :

In [1]:

```
1 from sklearn.datasets import make_moons
2 import matplotlib.pyplot as plt
3 import matplotlib
4 from matplotlib.colors import ListedColormap
5 import numpy as np
6 matplotlib.rcParams['figure.figsize'] = (6.0, 6.0) # pour avoir des figures d
7 np.random.seed(0)
8 X, y = make_moons(n_samples=1000, noise=0.1)
9 cm_bright = ListedColormap(['#FF0000', '#0000FF'])
10 plt.scatter(X[:,0], X[:,1], s=40, c=y, cmap=cm_bright)#cmap=plt.cm.PiYG)
```

Out[1]:

<matplotlib.collections.PathCollection at 0x118c77860>

In [2]:

```
1  #from :
2  #https://github.com/ardendertat/Applied-Deep-Learning-with-Keras/blob/master/
3
4  def plot_decision_boundary(func, X, y):
5      amin, bmin = X.min(axis=0) - 0.1
6      amax, bmax = X.max(axis=0) + 0.1
7      hticks = np.linspace(amin, amax, 101)
8      vticks = np.linspace(bmin, bmax, 101)
9      aa, bb = np.meshgrid(hticks, vticks)
10     ab = np.c_[aa.ravel(), bb.ravel()]
11     c = func(ab)
12     cc = c.reshape(aa.shape)
13     cm = plt.cm.RdBu
14     cm_bright = ListedColormap(['#FF0000', '#0000FF'])
15     fig, ax = plt.subplots()
16     contour = plt.contourf(aa, bb, cc, cmap=cm, alpha=0.8)
17     ax_c = fig.colorbar(contour)
18     ax_c.set_label("$P(y = 1)$")
19     ax_c.set_ticks([0, 0.25, 0.5, 0.75, 1])
20
21     plt.scatter(X[:, 0], X[:, 1], c=y, cmap=cm_bright)
22     plt.xlim(amin, amax)
23     plt.ylim(bmin, bmax)
24     plt.title("Decision Boundary")
```

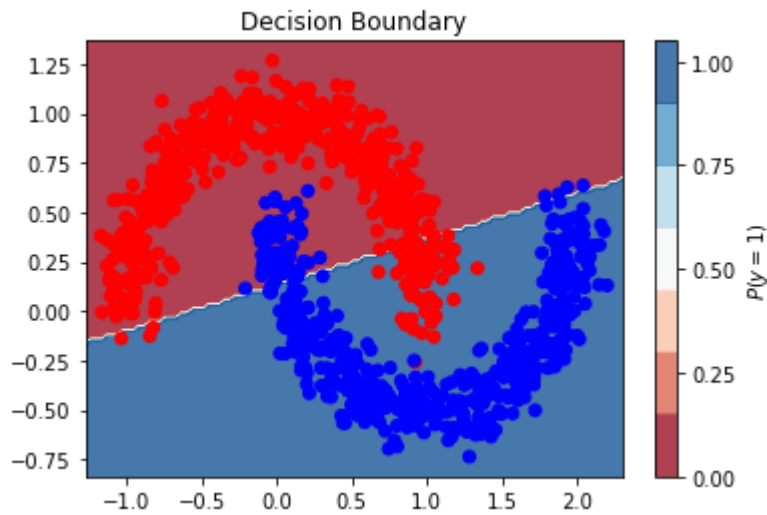
Appliquons, à présent, la logistic regression de sickit learn pour afficher la frontière de decision :

In [3]:

```

1  from sklearn import linear_model
2  clf = linear_model.LogisticRegression(solver='lbfgs')
3  clf.fit(X, y)
4  #Plot the decision boundary
5  plot_decision_boundary(lambda x: clf.predict(x), X, y)

```



Comme nous pouvons le constater les deux ensembles ne peuvent pas être séparés linéairement. Nous avons besoin de quelque chose de plus sophistiqué : les réseaux de neurones.

## Les réseaux de neurones

Les réseaux de neurones se composent des éléments suivants :

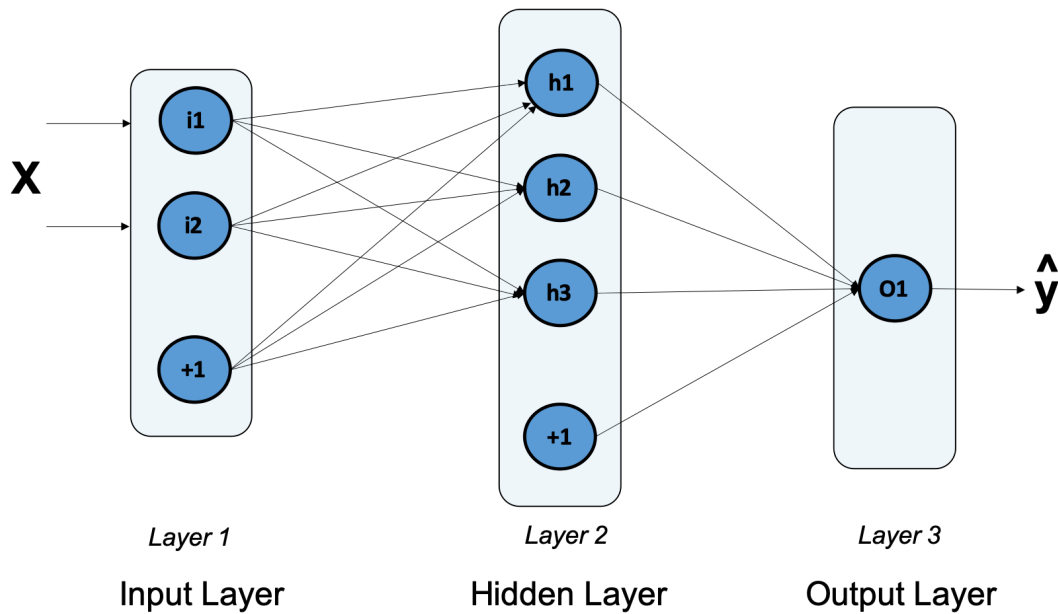
- Une couche d'entrée qui reçoit l'ensemble des caractéristiques (features), i.e. les variables prédictives.
- Un nombre arbitraire de couches cachées.
- Une couche de sortie,  $\hat{y}$ , qui contient la variable à prédire.
- Un ensemble de poids  $W$  qui vont être ajoutés aux valeurs des features et de biais  $b$  entre chaque couche
- Un choix de fonction d'activation pour chaque couche cachée,  $\sigma$ .

*Remarque* La couche de sortie doit avoir autant de neurones qu'il y a de sorties au problème de classification :

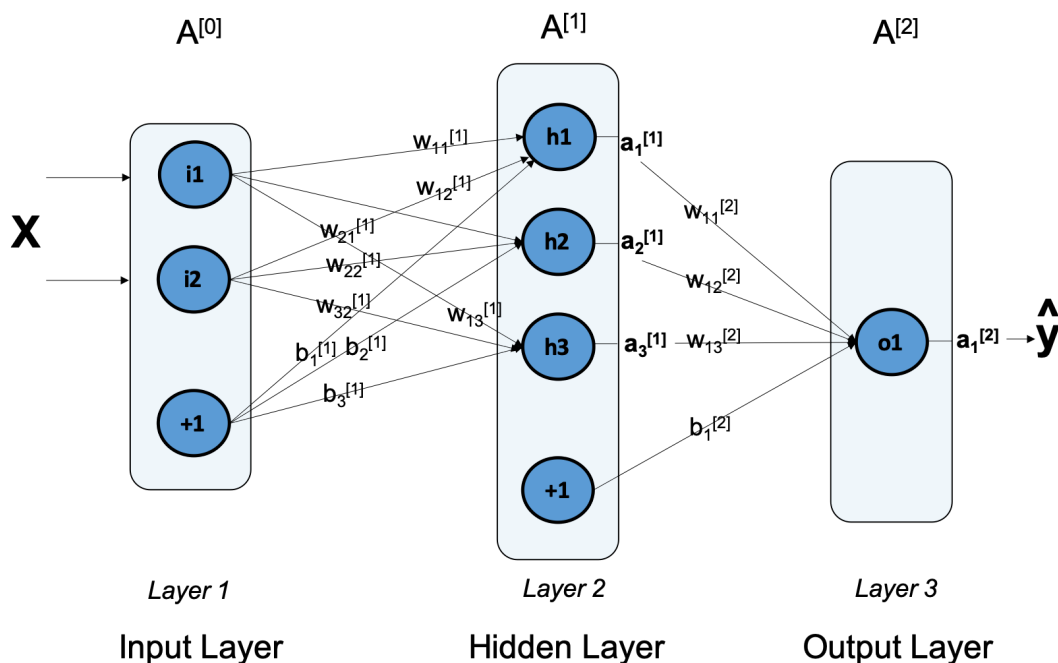
- *régression* : 1 seul neurone (C.f. notebook descente de gradient)
- *classification binaire* : 1 seul neurone avec une fonction d'activation qui sépare les deux classes.
- *classification multi-classe* : 1 neurone par classe et une fonction d'activation Softmax pour avoir la classe appropriée en fonction des probabilités de l'entrée appartenant à chaque classe.

La figure suivante illustre un exemple de réseau avec 3 couches :

- le layer 1 correspond au layer d'entrée (*input layer*), il reçoit l'ensemble des variables prédictives et est composé de 2 neurones. Le neurone avec +1 correspond au biais qui est ajouté.
- le layer 2 est appelé couche cachée (*hidden layer*), il possède 3 neurones et aussi un biais.
- le layer 3 correspond à la couche de sortie (*output layer*), la sortie de ce layer correspond à la prédiction.

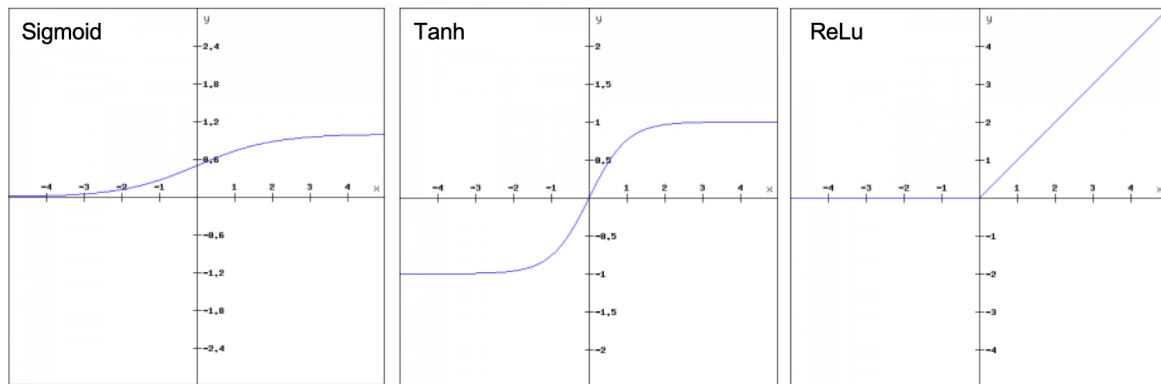


La figure suivante illustre le même réseau avec les poids affectés.



**Notations :**  $X$  correspond aux variables prédictives. Le poids est identifié de la manière suivante :  $w_{ij}^{[l]}$  où  $l$  correspond au niveau du layer cible,  $i$  correspond au numéro du nœud de la connection dans la couche  $l - 1$  et  $j$  correspond au numéro du nœud de la connection dans la couche  $l$ . Par exemple, le poids entre le nœud 1 dans le layer 1 et le nœud 2 dans le layer 2 est noté :  $w_{12}^{[2]}$ . Un biais est connecté à chaque nœud de la couche suivante. La notation est similaire :  $b_i^{[l]}$  où  $i$  est le numéro du nœud de la couche supérieure. La sortie d'un nœud est notée  $a_i^{[l]}$  où  $i$  correspond au numéro du nœud dans la couche  $l$ .  $\hat{y}$  correspond à la variable prédite.

### Choix de la fonction d'activation

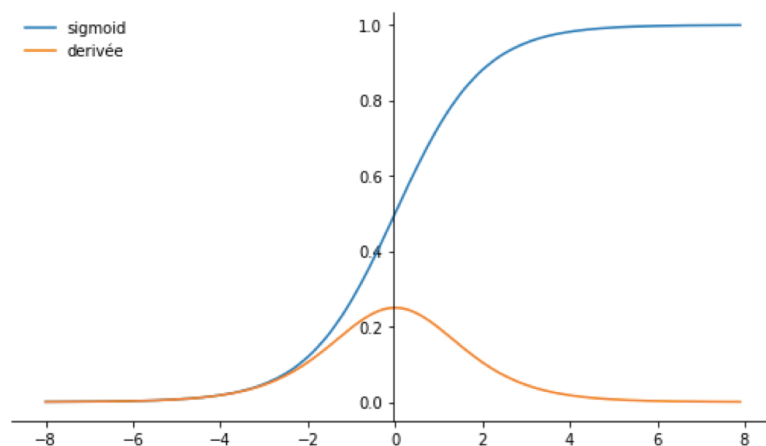


Il existe de très nombreuses fonctions d'activation qui peuvent être utilisées :

- Binary Step
- Sigmoid
- Tanh
- ReLU
- Leaky ReLU
- Softmax
- ...

Elles n'ont pas les mêmes propriétés.

Nous verrons, par la suite, que les réseaux de neurones utilisent la descente de gradient, le comportement de la dérivée des fonctions est donc important. Par exemple, nous avons vu que la sigmoid va transformer de grandes valeurs d'entrée dans des valeurs comprises entre 0 et 1. Cela veut dire qu'une modification importante de l'entrée entraînera une modification mineure de la sortie (C.f. notebook descente de gradient). Par conséquent la dérivée devient plus petite comme l'illustre l'image ci-dessous :



En fait, pour corriger les erreurs, les dérivées du réseau vont être propagées layer par layer de l'output layer à l'input layer. Le problème est que celles-ci sont multipliées entre chaque layer afin de connaître les valeurs de dérivées utiles pour l'input layer : le gradient décroît de façon exponentielle à mesure que nous nous propageons jusqu'aux couches initiales.

Pour choisir les fonctions d'activation, il faut considérer les propriétés principales suivantes :

- La disparition du gradient (*vanishing gradient*) : le problème intervient généralement dans des réseaux avec de très nombreux layer. Comme les descentes de gradient sont propagées dans tout le réseau, de trop petites valeurs de gradient (le gradient de la fonction de perte approche 0) indiquent que les poids des premiers layers ne seront pas mis à jour efficacement à chaque étape. Ceci entraîne donc une imprécision globale du réseau. Cela peut arriver si le réseau est composé de nombreuses couches avec une sigmoid.

- Disparition de neurones (*dead neuron*) : un neurone mort est un neurone qui, lors de l'apprentissage, ne s'active plus. Cela est lié au fait que les dérivées sont très petites ou nulles. Le neurone ne peut donc pas mettre à jour les poids. Les erreurs ne se propageant plus, ce neurone peut affecter les autres neurones du réseau. C'est, par exemple, le cas avec ReLu qui renvoie 0 quand l'entrée est inférieure ou égale à 0. Si chaque exemple donne une valeur négative, le neurone ne s'active pas et après la descente de gradient le neurone devient 0 donc ne sera plus utilisé. Le Leaky Relu permet de résoudre ce problème.
- Explosion du gradient (*Exploding gradient*) : le problème se pose lorsque des gradients d'erreur important s'accumulent et entraînent des mises à jour importantes des poids. Cela amène un réseau instable : les valeurs de mises à jour des poids peuvent être trop grandes et être remplacées par des NaN donc non utilisables (s'il n'y a pas d'erreurs d'exécution bien sûr !). Le problème est lié au type de descente de gradient utilisé (Batch vs mini-batch), au fait qu'il y a peut être trop de couches dans le réseau et bien sûr à certaines fonctions d'activation qui favorisent ce problème.
- Saturation de neurones (*Saturated neurons*) : le problème est lié au fait que les valeurs grandes (resp. petites) atteignent un plafond et qu'elles ne changent pas lors de la propagation dans le réseau. Ce problème est principalement lié aux fonctions sigmoid et tanh. En effet, sigmoid, pour toutes les valeurs supérieures à 1 va arriver sur un plateau et retournera toujours 1. Pour cela, ces deux fonctions d'activations sont assez déconseillées en deep learning (préférer LeRu ou Leaky Relu).

Pour avoir une idée du comportement des différentes fonctions d'activation et de leurs conséquences :

Une fois que tout est fixe, le réseau de neurones s'exécute alors en deux étapes :

- Forward Propagation
- Backward Propagation

## Forward Propagation

L'objectif de cette étape est de déterminer la valeur de sortie du réseau :  $\hat{\mathbf{y}}$ .

Comme nous avons vu dans le notebook descente de gradient, pour chaque neurone de la couche, nous effectuons une application affine en considérant les valeurs issues de la couche précédente (i.e.  $a$  représente le résultat de la fonction d'activation de la couche précédente) :

$$\mathbf{z}_i^{[l]} = \mathbf{w}_i^T \cdot \mathbf{a}^{[l-1]} + b_i \quad \mathbf{a}_i^{[l]} = \sigma^{[l]}(\mathbf{z}_i^{[l]})$$

Que nous pouvons donc généraliser en utilisant les matrices :

$$\begin{aligned} \mathbf{Z}^{[l]} &= \mathbf{W}^{[l]} \cdot \mathbf{A}^{[l-1]} + \mathbf{b}^{[l]} \\ \mathbf{A}^{[l]} &= \sigma^{[l]}(\mathbf{Z}^{[l]}) \end{aligned}$$

Si nous reprenons l'exemple de réseau précédent avec ReLu pour le hidden layer et sigmoid pour le layer de sortie nous avons donc :

$$\mathbf{A}^{[0]} = \mathbf{X}$$

où  $\mathbf{A}^{[0]} = \mathbf{X}$  la matrice contenant les exemples d'apprentissage.

$$\begin{aligned} \mathbf{Z}^{[1]} &= \mathbf{W}^{[1]} \cdot \mathbf{A}^{[0]} + \mathbf{b}^{[1]} \\ \mathbf{A}^{[1]} &= \text{ReLU}^{[1]}(\mathbf{Z}^{[1]}) \\ \mathbf{Z}^{[2]} &= \mathbf{W}^{[2]} \cdot \mathbf{A}^{[1]} + \mathbf{b}^{[2]} \\ \mathbf{A}^{[2]} &= \text{Sigmoid}^{[2]}(\mathbf{Z}^{[2]}) \end{aligned}$$

Finalement :  $\hat{y} = \mathbf{A}^{[2]}$

Le code ci-dessous illustre un exemple simple de la phase forward propagation pour notre exemple comportant deux variables prédictives. Les fonctions d'activations sont respectivement Relu pour le hidden layer et sigmoid pour le dernier layer.

In [4]:

```

1  import numpy as np
2
3  #fonctions d'activation
4  def sigmoid(Z):
5      return 1/(1+np.exp(-Z))
6
7  def relu(Z):
8      return np.maximum(0,Z)
9
10 def valueoutput(y_hat):
11     for i in range(len(y_hat)):
12         if y_hat[i]>0.5:
13             y_hat[i]=1
14         else:
15             y_hat[i]=0
16     return y_hat
17
18 # les donnees d'entree sous la forme d'une matrice (array en python)
19 X = np.array([0, 1],
20              [0, 0],
21              [1, 0],
22              [1, 1],
23              [1, 1]))
24 # les donnees de sortie sous la forme d'un vecteur
25 y = np.array([1],
26              [0],
27              [0],
28              [1],
29              [1]), dtype=float)
30
31
32 # initialisation des poids de manière aléatoire ainsi que des biais
33 inputSize = 2
34 hiddenSize = 3
35 outputSize = 1
36 W1=np.random.rand(2, 3)
37 W2=np.random.rand(3, 1)
38 b1 = np.random.rand(3)
39 b2 = np.random.rand(1)
40 print ("Les données d'entrées : \n",X)
41 print ('Les valeurs de poids et de biais initialisees aléatoirement : \n')
42 print ('\t(layer input vers layer 1) : W1 \n',W1,'\n')
43 print ('\t(layer input vers layer 1) : b1\n',b1,'\n')
44 print ("\t(layer 1 vers layer 2) : W2\n",W2,'\n')
45 print ("\t(layer 1 vers layer 2) : b2\n",b2,'\n')
46 print ("Etape 1 : ")
47 print ("\n A0=X\n")
48
49 A0=X
50 Z1 = np.dot(A0,W1)+b1
51 print ("\n Z1 = W1.A0 + b1 \n",Z1,'\n')
52 A1 = relu(Z1)
53 print ('\nA1 = relu(Z1)\n',A1,'\n')
54 Z2 = np.dot(A1,W2)+b2
55 print ("\n Z2 = W2.A1 + b2 \n",Z2,'\n')
56 A2 = sigmoid(Z2)
57 print ('\nA2 = sigmoid(Z2)\n',A2,'\n')
58 y_hat = sigmoid(Z2)
59 print ('yhat\n',y_hat)

```



```

60
61
62
63     print ("Les données d'entrées : \n",X)
64     print ("Les sorties predites : \n", str(valueoutput(y_hat)))
65
66     print ("Les sorties reelles attendues : \n", str(y))

```

Les données d'entrées :

```

[[0 1]
 [0 0]
 [1 0]
 [1 1]
 [1 1]]

```

Les valeurs de poids et de biais initialisees aléatoirement :

```

(layer input vers layer 1) : W1
[[0.90496764 0.66934312 0.61669425]
 [0.70675322 0.21538845 0.58636595]]

```

```

(layer input vers layer 1) : b1
[0.55441685 0.50237972 0.86147085]

```

```

(layer 1 vers layer 2) : W2
[[0.68936531]
 [0.25632519]
 [0.84027211]]

```

```

(layer 1 vers layer 2) : b2
[0.18734206]

```

Etape 1 :

A0=X

```

Z1 = W1.A0 + b1
[[1.26117007 0.71776817 1.4478368 ]
 [0.55441685 0.50237972 0.86147085]
 [1.45938449 1.17172285 1.4781651 ]
 [2.16613771 1.38711129 2.06453105]
 [2.16613771 1.38711129 2.06453105]]

```

```

A1 = relu(Z1)
[[1.26117007 0.71776817 1.4478368 ]
 [0.55441685 0.50237972 0.86147085]
 [1.45938449 1.17172285 1.4781651 ]
 [2.16613771 1.38711129 2.06453105]
 [2.16613771 1.38711129 2.06453105]]

```

```

Z2 = W2.A1 + b2
[[2.45730789]
 [1.4221803 ]
 [2.73579409]
 [3.77092168]
 [3.77092168]]

```

A2 = sigmoid(Z2)

```

rrr 921094221

```

```
[[0.92109422]
 [0.80567999]
 [0.93910602]
 [0.97748765]
 [0.97748765]]
```

```
yhat
[[0.92109422]
 [0.80567999]
 [0.93910602]
 [0.97748765]
 [0.97748765]]
```

Les données d'entrées :

```
[[0 1]
 [0 0]
 [1 0]
 [1 1]
 [1 1]]
```

Les sorties predites :

```
[[1.]
 [1.]
 [1.]
 [1.]
 [1.]]
```

Les sorties reelles attendues :

```
[[1.]
 [0.]
 [0.]
 [1.]
 [1.]]
```

Comme nous pouvons le constater il y a des erreurs dans les sorties prédites. C'est là qu'intervient la seconde phase.

## Backward Propagation

L'objectif de la Backward Propagation est tout d'abord d'évaluer la différence entre la valeur prédite et la valeur réelle.

*Etape 1 : (calcul du coût)*

Nous avons vu dans le notebook de la descente de gradient que la différence entre la valeur obtenue dans l'étape précédente et la valeur réelle correspond au coût. Plus la différence est élevée, plus le coût sera élevé. Pour minimiser ce coût, il faut trouver les valeurs de poids et de biais pour lesquelles la fonction de coût renvoie la plus petite valeur possible. Plus le coût est faible, plus les prévisions sont exactes. Nous retrouvons donc le problème rencontré pour la descente de gradient.

Précédemment nous avons vu que la cross entropy, comme fonction de coût, était bien adaptée à notre problème de classification binaire, donc :

$$C(\hat{y}, y) = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y})$$

L'objectif, à présent, est de propager cette erreur dans tout le réseau pour mettre à jour les différents poids.

*Etape 2 : (backpropagation)*

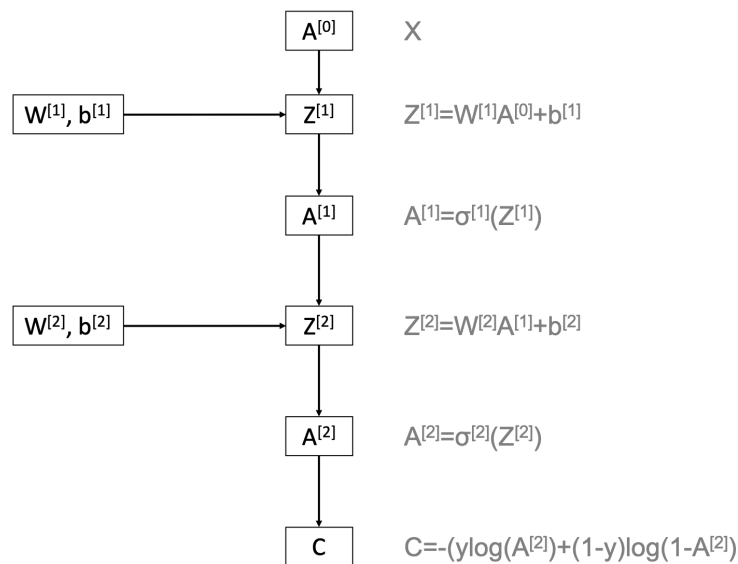
**Comprendre ce qui est derrière**

Nous avons vu précédemment, lors de la phase de forward, que l'exécution était de la forme :

$$\begin{aligned}
 \mathbf{A}^{[0]} &= \mathbf{X} \\
 \mathbf{Z}^{[1]} &= \mathbf{W}^{[1]} \cdot \mathbf{A}^{[0]} + \mathbf{b}^{[1]} \\
 \mathbf{A}^{[1]} &= \sigma^{[1]}(\mathbf{Z}^{[1]}) \\
 \mathbf{Z}^{[2]} &= \mathbf{W}^{[2]} \cdot \mathbf{A}^{[1]} + \mathbf{b}^{[2]} \\
 \mathbf{A}^{[2]} &= \sigma^{[2]}(\mathbf{Z}^{[2]}) \\
 &\vdots \\
 \mathbf{Z}^{[L]} &= \mathbf{W}^{[L]} \cdot \mathbf{A}^{[L-1]} + \mathbf{b}^{[L]} \\
 \mathbf{A}^{[L]} &= \sigma^{[L]}(\mathbf{Z}^{[L]}) = \hat{y}
 \end{aligned}$$

où  $L$  est le output layer.

La figure suivante illustre les étapes jusqu'à la fonction de coût pour notre réseau exemple :



L'objectif de la backward propagation est de reporter, dans le réseau, l'ensemble des modifications à apporter aux poids entre les couches. Pour cela il faut repartir en sens inverse pour calculer les dérivées partielles du coût. Elle repose sur la règle de dérivation en chaîne (*chain rule*) qui est une formule qui explicite la dérivée d'une fonction composée pour deux fonctions dérivables :

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

Lorsque l'on regarde la fin du réseau, nous constatons que  $\mathbf{C}$  est une fonction qui dépend de  $\mathbf{A}^{[2]}$ , que  $\mathbf{A}^{[2]}$  dépend, elle-même, d'une fonction  $\mathbf{Z}^{[2]}$  et que finalement  $\mathbf{Z}^{[2]}$  dépend de  $\mathbf{W}^{[2]}$  et de  $\mathbf{b}^{[2]}$ .

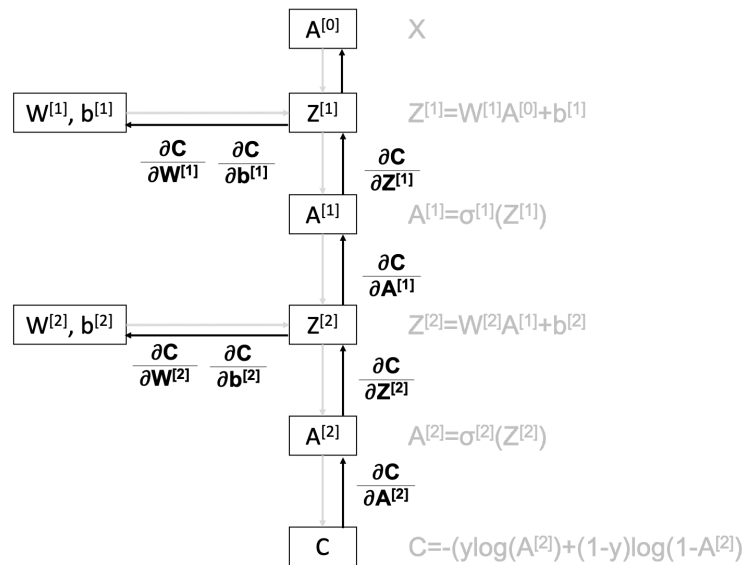
$$\begin{aligned}
 \frac{\partial \mathbf{C}}{\partial \mathbf{W}^{[2]}} &= \frac{\partial \mathbf{C}}{\partial \mathbf{A}^{[2]}} \cdot \frac{\partial \mathbf{A}^{[2]}}{\partial \mathbf{Z}^{[2]}} \cdot \frac{\partial \mathbf{Z}^{[2]}}{\partial \mathbf{W}^{[2]}} \\
 \frac{\partial \mathbf{C}}{\partial \mathbf{b}^{[2]}} &= \frac{\partial \mathbf{C}}{\partial \mathbf{A}^{[2]}} \cdot \frac{\partial \mathbf{A}^{[2]}}{\partial \mathbf{Z}^{[2]}} \cdot \frac{\partial \mathbf{Z}^{[2]}}{\partial \mathbf{b}^{[2]}}
 \end{aligned}$$

De la même manière, pour avoir la dérivée partielle de  $\mathbf{C}$  par rapport à  $\mathbf{W}^{[1]}$  et  $\mathbf{b}^{[1]}$ , nous voyons, sur la figure, que  $\mathbf{Z}^{[2]}$  est une fonction qui dépend de  $\mathbf{A}^{[1]}$ , qui, elle-même, dépend de  $\mathbf{Z}^{[1]}$  et que finalement  $\mathbf{Z}^{[1]}$  dépend de  $\mathbf{W}^{[1]}$  et  $\mathbf{b}^{[1]}$ .

$$\frac{\partial \mathbf{C}}{\partial \mathbf{W}^{[1]}} = \frac{\partial \mathbf{C}}{\partial \mathbf{A}^{[2]}} \cdot \frac{\partial \mathbf{A}^{[2]}}{\partial \mathbf{Z}^{[2]}} \cdot \frac{\partial \mathbf{Z}^{[2]}}{\partial \mathbf{A}^{[1]}} \cdot \frac{\partial \mathbf{A}^{[1]}}{\partial \mathbf{Z}^{[1]}} \cdot \frac{\partial \mathbf{Z}^{[1]}}{\partial \mathbf{W}^{[1]}}$$

$$\frac{\partial C}{\partial \mathbf{b}^{[1]}} = \frac{\partial C}{\partial \mathbf{A}^{[2]}} \cdot \frac{\partial \mathbf{A}^{[2]}}{\partial \mathbf{Z}^{[2]}} \cdot \frac{\partial \mathbf{Z}^{[2]}}{\partial \mathbf{A}^{[1]}} \cdot \frac{\partial \mathbf{A}^{[1]}}{\partial \mathbf{Z}^{[1]}} \cdot \frac{\partial \mathbf{Z}^{[1]}}{\partial \mathbf{b}^{[1]}}$$

Les différentes étapes sont résumées sur la figure suivante :



Pour résumer, les équations pour calculer la dérivée partielle de la fonction de coût en fonction des poids et des biais d'une couche  $l$  sont :

$$\frac{\partial C}{\partial W^{[l]}} = \frac{\partial C}{\partial Z^{[l]}} \cdot \frac{\partial Z^{[l]}}{\partial W^{[l]}}$$

$$\frac{\partial C}{\partial b^{[l]}} = \frac{\partial C}{\partial Z^{[l]}} \cdot \frac{\partial Z^{[l]}}{\partial b^{[l]}}$$

Donc, pour obtenir les dérivées partielles de  $C$  par rapport à  $W^{[l]}$  et  $b^{[l]}$ , nous devons calculer :

$$\frac{\partial C}{\partial A^{[l]}}, \frac{\partial C}{\partial Z^{[l]}}, \frac{\partial C}{\partial Z^{[l]}} \frac{\partial Z^{[l]}}{\partial W^{[l]}}, \frac{\partial Z^{[l]}}{\partial b^{[l]}}$$

Par la suite, et par simplification, nous considérons que les deux fonctions d'activation dans notre réseau sont Relu (pour  $A^{[1]}$ ) et sigmoid (pour  $A^{[2]}$ ). Le principe est le même quelques soient les fonctions, il suffit juste de connaître la dérivée des fonctions d'activation.

**Pour la dérivée partielle de  $C$  par rapport à  $A^{[L]}$  :**

$$\frac{\partial C}{\partial A^{[L]}} = \frac{\partial (-y \log(A^{[L]}) - (1-y) \log(1-A^{[L]}))}{\partial A^{[L]}}$$

Nous savons que la dérivée d'une fonction  $\log(x)$  est :

$$\frac{\partial \log(x)}{\partial x} = \frac{1}{x}$$

Pour la partie gauche  $-y \log(A^{[L]})$  nous avons donc comme dérivée :

$$\frac{-y}{A^{[L]}}$$

Pour la partie droite  $-(1-y) \log(1-A^{[L]})$ , il faut juste appliquer la formule de la dérivée d'une fonction :

$$\frac{\partial \log(g(x))}{\partial x} = \frac{1}{g(x)} g'(x)$$

comme la dérivée de  $1 - A^{[L]}$  est  $-1$  nous avons au final :

$$\begin{aligned}\frac{\partial C}{\partial A^{[L]}} &= \frac{-y}{A^{[L]}} - (-) \frac{(1-y)}{(1-A^{[L]})} \\ &= \left( \frac{-y}{A^{[L]}} + \frac{(1-y)}{(1-A^{[L]})} \right)\end{aligned}$$

Donc :

$$\boxed{\frac{\partial C}{\partial A^{[L]}} = \left( \frac{-y}{A^{[L]}} + \frac{(1-y)}{(1-A^{[L]})} \right)}$$

**Considérons, à présent la dérivée partielle de C par rapport à  $Z^{[L]}$  :**

$$\frac{\partial C}{\partial Z^{[L]}}$$

En utilisant la chaîne de dérivation :

$$\frac{\partial C}{\partial Z^{[L]}} = \frac{\partial C}{\partial A^{[L]}} \cdot \frac{\partial A^{[L]}}{\partial Z^{[L]}} = \frac{\partial C}{\partial A^{[L]}} * \sigma'^{[L]}(Z^{[L]})$$

$\sigma'^{[L]}(Z^{[L]})$  correspond simplement à la dérivée de la sigmoid. Nous avons vu dans le notebook sur la descente de gradient, que cette dérivée est :

$$\frac{\partial A^{[L]}}{\partial Z^{[L]}} = \text{sigmoid}(Z^{[L]})(1 - \text{sigmoid}(Z^{[L]})) = A^{[L]}(1 - A^{[L]})$$

Donc :

$$\frac{\partial C}{\partial A^{[L]}} \cdot \frac{\partial A^{[L]}}{\partial Z^{[L]}} = \left( \frac{-y}{A^{[L]}} + \frac{(1-y)}{(1-A^{[L]})} \right) A^{[L]}(1 - A^{[L]})$$

Nous pouvons multiplier par  $(1 - A^{[L]})$  et  $(A^{[L]})$  pour simplifier :

$$\begin{aligned}&= \left( \frac{-y(1 - A^{[L]})}{A^{[L]}(1 - A^{[L]})} + \frac{A^{[L]}(1 - y)}{A^{[L]}(1 - A^{[L]})} \right) A^{[L]}(1 - A^{[L]}) \\ &= \left( \frac{-y(1 - A^{[L]}) + A^{[L]}(1 - y)}{A^{[L]}(1 - A^{[L]})} \right) A^{[L]}(1 - A^{[L]})\end{aligned}$$

en supprimant  $A^{[L]}(1 - A^{[L]})$  nous avons :

$$\begin{aligned}&= (-y(1 - A^{[L]}) + A^{[L]}(1 - y)) \\ &= -y + yA^{[L]} + A^{[L]} - A^{[L]}y \\ &= -y + A^{[L]}\end{aligned}$$

Donc :

$$\boxed{\frac{\partial C}{\partial Z^{[L]}} = A^{[L]} - y}$$

**Dérivée partielle de C par rapport à  $Z^{[l]}$  :**

Nous souhaitons, à présent, obtenir la dérivée partielle de C par rapport à un niveau l, i.e.  $Z^{[l]}$ . Le principe étant que si l'on connaît  $Z^{[L]}$ , il est possible de déduire  $Z^{[L-1]}$ ,  $Z^{[L-2]}$ , ...

Nous savons, en appliquant la chaîne de dérivation, qu'il est possible de calculer la dérivée de  $\frac{\partial C}{\partial Z^{[l]}}$  par :

$$\frac{\partial C}{\partial Z^{[l]}} = \frac{\partial C}{\partial Z^{[l+1]}} \cdot \frac{\partial Z^{[l+1]}}{\partial A^{[l]}} \cdot \frac{\partial A^{[l]}}{\partial Z^{[l]}}$$

Comme :

$$\begin{aligned} Z^{[l+1]} &= W^{[l+1]} \cdot A^{[l]} + b^{[l+1]} \\ \frac{\partial Z^{[l+1]}}{\partial A^{[l]}} &= \frac{\partial (W^{[l+1]} \cdot A^{[l]} + b^{[l+1]})}{\partial A^{[l]}} \\ &= W^{[l+1]} \end{aligned}$$

Nous avons également :

$$\frac{\partial A^{[l]}}{\partial Z^{[l]}} = \sigma'^{[l]}(Z^{[l]})$$

Donc :

$$\boxed{\frac{\partial C}{\partial Z^{[l]}} = (W^{[l+1]})^T \cdot \frac{\partial C}{\partial Z^{[l+1]}} * \sigma'^{[l]}(Z^{[l]})}$$

*Dérivée partielle de  $Z^{[l]}$  par rapport à  $W^{[l]}$  :*

Dans un premier temps nous calculons la dérivée partielle de  $Z^{[l]}$  par rapport à  $W^{[l]}$  pour, par la suite déterminer la dérivée partielle de  $C$  par rapport à  $W^{[l]}$ .

Comme :

$$\begin{aligned} Z^{[l]} &= W^{[l]} \cdot A^{[l-1]} + b^{[l]} \\ \frac{\partial Z^{[l]}}{\partial W^{[l]}} &= \frac{\partial (W^{[l]} \cdot A^{[l-1]} + b^{[l]})}{\partial W^{[l]}} \\ &= A^{[l-1]} \end{aligned}$$

**Dérivée partielle de  $C$  par rapport à  $W^{[l]}$  :**

A partir du résultat précédent, nous avons :

$$\boxed{\frac{\partial C}{\partial W^{[l]}} = \frac{\partial C}{\partial Z^{[l]}} \cdot A^{[l-1]T}}$$

*Dérivée partielle de  $Z^{[l]}$  par rapport à  $b^{[l]}$  :*

Comme précédemment, nous calculons la dérivée partielle de  $Z^{[l]}$  par rapport à  $b^{[l]}$  pour, par la suite déterminer la dérivée partielle de  $C$  par rapport à  $b^{[l]}$ .

Comme :

$$\begin{aligned} Z^{[l]} &= W^{[l]} \cdot A^{[l-1]} + b^{[l]} \\ \frac{\partial Z^{[l]}}{\partial b^{[l]}} &= \frac{\partial (W^{[l]} \cdot A^{[l-1]} + b^{[l]})}{\partial b^{[l]}} \\ &= 1 \end{aligned}$$

**Dérivée partielle de  $C$  par rapport à  $b^{[l]}$  :**

A partir du résultat précédent. nous avons :

À partir du résultat précédent, nous avons :

$$\frac{\partial C}{\partial \mathbf{b}^{[l]}} = \frac{\partial C}{\partial \mathbf{z}^{[l]}}$$

### Pour résumer

À présent, pour le dernier layer  $L$ , nous sommes capable de calculer la dérivée partielle du coût par rapport à  $\mathbf{A}^{[L]}$ ,  $\mathbf{z}^{[L]}$ ,  $\mathbf{W}^{[L]}$  et  $\mathbf{b}^{[L]}$  :

$\frac{\partial C}{\partial \mathbf{A}^{[L]}}$	$\left( \frac{-y}{\mathbf{A}^{[L]}} + \frac{(1-y)}{(1-\mathbf{A}^{[L]})} \right)$
$\frac{\partial C}{\partial \mathbf{z}^{[L]}}$	$(\mathbf{A}^{[L]} - y)$
$\frac{\partial C}{\partial \mathbf{W}^{[L]}}$	$\frac{\partial C}{\partial \mathbf{z}^{[L]}} \cdot (\mathbf{A}^{[L-1]})^T$
$\frac{\partial C}{\partial \mathbf{b}^{[L]}}$	$\frac{\partial C}{\partial \mathbf{z}^{[L]}}$

Pour n'importe quel layer  $l$ , nous avons :

$\frac{\partial C}{\partial \mathbf{z}^{[l]}}$	$(\mathbf{W}^{[l+1]})^T \cdot \frac{\partial C}{\partial \mathbf{z}^{[l+1]}} * \sigma'^{[l]}(\mathbf{z}^{[l]})$
$\frac{\partial C}{\partial \mathbf{W}^{[l]}}$	$\frac{\partial C}{\partial \mathbf{z}^{[l]}} \cdot \mathbf{A}^{[l-1]T}$
$\frac{\partial C}{\partial \mathbf{b}^{[l]}}$	$\frac{\partial C}{\partial \mathbf{z}^{[l]}}$

### La descente de gradient

Attention, parfois, la backward propagation est considérée comme la descente de gradient. Ce n'est pas le cas. Elle a pour seul objectif de calculer les gradients pour les opérations à chaque niveau. La descente de gradient intervient après, elle permet de pouvoir mettre automatiquement les poids des différents layer en appliquant justement les gradients obtenus dans l'étape précédente.

La descente de gradient se fait comme dans le notebook : il faut boucler jusqu'au premier layer pour appliquer la formule du gradient à l'aide des dérivées calculées précédemment :

**For  $l$  in enumerate (dernier\_layer,1) {**

$$\mathbf{W}^{[l]} = \mathbf{W}^{[l]} - \eta \frac{\partial C}{\partial \mathbf{W}^{[l]}}$$

$$\mathbf{b}^{[l]} = \mathbf{b}^{[l]} - \eta \frac{\partial C}{\partial \mathbf{b}^{[l]}}$$

**}**

Remarque : comme nous l'avons vu lors des dérivations, il est nécessaire de sauvegarder  $\frac{\partial C}{\partial \mathbf{W}^{[l]}}$  et  $\frac{\partial C}{\partial \mathbf{b}^{[l]}}$  pour pouvoir les réutiliser lors de la descente de gradient.

### Cas de la classification multi-classes

Jusqu'à présent nous avons vu comment faire de la classification binaire, i.e. la fonction d'activation est une sigmoid. Pour faire de la classification multi-classe, il faut utiliser la fonction d'activation softmax. Elle attribue des probabilités à chaque classe d'un problème à plusieurs classes et la somme de ces probabilités doit être égale à 1. Formellement softmax, prend en entrée un vecteur de  $C$ -dimensions (le nombre de classes possibles)  $\mathbf{z}$  et retourne un autre vecteur de  $C$ -dimensions  $\mathbf{a}$  de valeurs réelles comprises entre 0 et 1.

Pour  $i = 1 \dots C$  :

$$a_i = \frac{e^{z_i}}{\sum_{k=1}^C e^{z_k}}$$

avec  $\sum_{i=1}^C a_i = 1$

où  $C$  correspond au nombre de classes.

In [5]:

```
1 def softmax(z):
2     expz = np.exp(z)
3     return expz / expz.sum(axis=0, keepdims=True)
4
5 nums = np.array([4, 5, 6])
6 print(softmax(nums))
7 print("la somme des probabilités donne 1")
```

```
[0.09003057 0.24472847 0.66524096]
la somme des probabilités donne 1
```

Cependant, cette fonction n'est pas très stable : elle génère souvent des nan pour des grands nombres par exemple.

In [6]:

```
1 nums = np.array([4000, 5000, 6000])
2 print(softmax(nums))
```

```
[nan nan nan]
```

```
/Users/pascalponcelet/Desktop/Sicki-learn/Tools/tools/lib/python3.6/site-packages/ipykernel_launcher.py:2: RuntimeWarning: overflow encountered in exp
```

```
/Users/pascalponcelet/Desktop/Sicki-learn/Tools/tools/lib/python3.6/site-packages/ipykernel_launcher.py:3: RuntimeWarning: invalid value encountered in true_divide
```

This is separate from the ipykernel package so we can avoid doing imports until

Aussi il est fréquent de multiplier le numérateur par une constante : généralement  $-max(z)$  :

$$a_i = \frac{e^{z_i - \max(z)}}{\sum_{k=1}^C e^{z_k - \max(z)}}$$



In [7]:

```

1  def softmax(z):
2      expz = np.exp(z - np.max(z))
3      return expz / expz.sum(axis=0, keepdims=True)
4
5  nums = np.array([4, 5, 6])
6  print(softmax(nums))
7  print("la somme des probabilités donne 1")
8  nums = np.array([4000, 5000, 6000])
9  print(softmax(nums))

```

```

[0.09003057 0.24472847 0.66524096]
la somme des probabilités donne 1
[0. 0. 1.]

```

### Dérivée de la fonction softmax

Elle est basée sur le fait de considérer la ré-écriture suivante :

$$g(x) = e^{z_i}$$

et

$$h(x) = \sum_{k=1}^C e^{z_k}$$

Nous savons que la dérivée d'une fonction  $f(x) = \frac{g(x)}{h(x)}$  est :

$$f'(x) = \frac{g'(x)h(x) - h'(x)g(x)}{h(x)^2}$$

Par simplification, nous notons :

$$\sum_C = \sum_{k=1}^C e^{z_k}$$

Pour  $i = 1 \dots C$ , nous avons :

$$a_i = \frac{e^{z_i}}{\sum_C}$$

La dérivée  $\frac{\partial a_i}{\partial z_j}$  de la sortie de softmax  $\mathbf{a}$  par rapport à  $\mathbf{z}$  :

- Si  $i = j$ ,

$$\frac{\partial a_i}{\partial z_i} = \frac{\partial(\frac{e^{z_i}}{\sum_C})}{\partial z_i} = \frac{e^{z_i} \sum_C - e^{z_i} e^{z_i}}{\sum_C^2} = \frac{e^{z_i}}{\sum_C} \frac{\sum_C - e^{z_i}}{\sum_C} = \frac{e^{z_i}}{\sum_C} (1 - \frac{e^{z_i}}{\sum_C}) = a_i(1 - a_i)$$

Remarque : il s'agit du même résultat que pour la dérivée de la sigmoid.

- Si  $i \neq j$ ,

$$\frac{\partial a_i}{\partial z_j} = \frac{\partial(\frac{e^{z_i}}{\sum_C})}{\partial z_j} = \frac{0 - e^{z_i} e^{z_j}}{\sum_C^2} = \frac{e^{z_i}}{\sum_C} \frac{e^{z_j}}{\sum_C} = -a_i a_j$$

### Dérivée par rapport à la fonction de coût

En appliquant le même principe que précédemment, on trouve :

$$\frac{\partial C}{\partial \mathbf{Z}^{[L]}} = \mathbf{A}^{[L]} - \mathbf{y}$$

Donc toutes les dérivées précédentes pour **W** et **b** sont également similaires.

Il est par contre nécessaire de redéfinir la fonction de prédiction. Généralement lorsqu'il y a plusieurs classes, il convient de transformer le **y** initial en utilisant la fonction *OneHotEncoder*. Parfois, au préalable, il est indispensable de transformer les labels s'il s'agit d'attributs catégoriels en nombre via *LabelEncoder*. (Cf. notebook ingénierie des données).

In [8]:

```

1  from sklearn.preprocessing import LabelEncoder
2  from sklearn.preprocessing import OneHotEncoder
3  import numpy as np
4
5  # Transformation du label
6  labelencoder = LabelEncoder()
7  exemple=np.array(["positif","negatif","positif","negatif"])
8  print ("Exemple ",exemple)
9  label_encoded=labelencoder.fit_transform(exemple)
10 print ("labels encodés :",label_encoded)
11
12
13 # Encodage via OneHotEncoder
14 onehot_encoder = OneHotEncoder(sparse=False,categories='auto')
15 #il est indispensable de faire un reshape
16
17 integer_encoded = label_encoded.reshape(len(label_encoded), 1)
18 final = onehot_encoder.fit_transform(integer_encoded)
19 print ("Encodage final :\n", final)

```

```

Exemple ['positif' 'negatif' 'positif' 'negatif']
labels encodés : [1 0 1 0]
Encodage final :
[[0. 1.]
 [1. 0.]
 [0. 1.]
 [1. 0.]]

```

Dans le cas de la prédiction, l'objectif est de retourner la classe qui a la plus forte probabilité à la sortie de softmax. Pour cela nous pouvons utiliser la fonction *argmax*.

In [9]:

```

1  exemple = np.array ([[0,2,25,4],
2                        [1,11,8,10],
3                        [200,7,5,10]])
4  print (exemple)
5

```

```

[[ 0  2 25  4]
 [ 1 11  8 10]
 [200  7  5 10]]

```

```

1  Position des différents éléments dans la matrice
2                        #0    #1    #2    #3
3  exemple = np.array ([[0,    2,   25,   4],    #0

```

```

4      [1, 11, 8, 10], #1
5      [200, 7, 5, 10]] #2
6  )

```

In [10]:

```

1  print ("\nPosition du plus grand élément : ", np.argmax(exemple))
2  print ("200 est le 8 ième élément (on compte à partir de 0)\n")
3
4  print ("Axis = 0 : la fonction cherche la valeur maximale sur les colonne de
5  print ("\nIndices de l'élément max en considérant les colonnes : ", np.argmax(e
6  print ("colonne 0 : 200 est le plus grand de la colonne (retourne ligne 2)")
7  print ("colonne 1 : 11 est le plus grand de la colonne (retourne ligne 1)")
8  print ("colonne 2 : 25 est le plus grand de la colonne (retourne ligne 0)")
9  print ("colonne 3 : 10 est le plus grand de la colonne (retourne ligne 1)\n")
10
11 print ("Axis = 1 : la fonction cherche la valeur maximale sur les lignes de l
12 print ("\nIndices of Max element : ", np.argmax(exemple, axis = 1))
13 print ("ligne 0 : 25 est le plus grand de la colonne (retourne colonne 2)")
14 print ("ligne 1 : 11 est le plus grand de la colonne (retourne colonne 1)")
15 print ("ligne 2 : 200 est le plus grand de la colonne (retourne colonne 0)")

```

Position du plus grand élément : 8  
 200 est le 8 ième élément (on compte à partir de 0)

Axis = 0 : la fonction cherche la valeur maximale sur les colonne de 1  
 a matrice

Indices de l'élément max en considérant les colonnes : [2 1 0 1]  
 colonne 0 : 200 est le plus grand de la colonne (retourne ligne 2)  
 colonne 1 : 11 est le plus grand de la colonne (retourne ligne 1)  
 colonne 2 : 25 est le plus grand de la colonne (retourne ligne 0)  
 colonne 3 : 10 est le plus grand de la colonne (retourne ligne 1)

Axis = 1 : la fonction cherche la valeur maximale sur les lignes de la  
 matrice

Indices of Max element : [2 1 0]  
 ligne 0 : 25 est le plus grand de la colonne (retourne colonne 2)  
 ligne 1 : 11 est le plus grand de la colonne (retourne colonne 1)  
 ligne 2 : 200 est le plus grand de la colonne (retourne colonne 0)

## Implémentations

### Implémentation sous la forme de fonctions

Récupération des librairies utiles. Dans la mesure du possible il est préférable de les mettre au tout début.

In [11]:

```

1  ▼  # importations utiles
2      import numpy as np
3      from sklearn.datasets import make_moons
4      import matplotlib.pyplot as plt
5      import matplotlib
6      from sklearn.model_selection import train_test_split
7      from sklearn import linear_model
8      from matplotlib.legend_handler import HandlerLine2D
9      import keras
10     from keras.models import Sequential
11     from keras.layers import Dense
12     from keras.utils import np_utils
13     from keras import regularizers
14     from keras.optimizers import SGD
15     from sklearn.metrics import accuracy_score
16
17     matplotlib.rcParams['figure.figsize'] = (6.0, 6.0) # pour avoir des figures d

```

Using TensorFlow backend.

Dans un premier temps il est nécessaire de construire le réseau en indiquant le nombre de layers, de neurones par layer et les fonctions d'activation. Nous stockons ici cette information sous la forme d'un dictionnaire python.

In [12]:

```

1  ▼  #reseau de l'exemple du notebook
2      # il contient trois layer : input - hidden - output
3      #   par défaut on ne spécifie pas la couche d'entrée (couche 0) qui contient
4      #   hidden layer : elle reçoit en entrée les variables prédictives X avec de
5      #   (input_dim=2) elle contient 5 neurones, la fonction d'activation
6      #   output layer : les dimensions d'entrée = 5 (il s'agit de la sortie de l'
7      #   elle donne le résultat. Ici la fonction d'activation est
8      #   classification binaire
9      #
10
11  ▼  layers = [
12      {"input_dim": 2, "output_dim": 3, "activation": "relu"},
13      {"input_dim": 3, "output_dim": 1, "activation": "sigmoid"}
14  ]
15
16

```

La fonction suivante permet de créer le réseau à partir du dictionnaire passé en paramètre. Elle initialise avec un nombre aléatoire les poids et le biais. Les différents paramètres du réseau (poids) et biais sont stockés sous la forme d'un dictionnaire python avec une clé qui identifie le layer auquel il appartient.

In [13]:

```

1  def init_layers(layers):
2
3      seed=30
4      np.random.seed(seed)
5
6      # nombre de layers dans le réseau
7      number_of_layers = len(layers)
8      # pour stocker les différentes valeurs des paramètres
9
10     paramameters = {}
11
12     # Pour toutes les couches du réseau
13     for idx, layer in enumerate(layers):
14         # Par simplification on commence la numérotation du layer à 1
15         # correspond aux données d'entrées (input), i.e. les variables prédic
16         layer_idx = idx + 1
17
18         layer_input_size = layer["input_dim"]
19         layer_output_size = layer["output_dim"]
20
21         # Initialisation des valeurs de la matrice W et du vecteur b
22         # pour les différentes couches
23
24         paramameters['W' + str(layer_idx)] = np.random.randn(
25             layer_output_size, layer_input_size) * 0.1
26         paramameters['b' + str(layer_idx)] = np.random.randn(
27             layer_output_size, 1) * 0.1
28
29     return paramameters

```

Définition des différentes fonctions d'activation (Relu, Sigmoid, Tanh) ainsi que les dérivées qui seront utilisées par la suite.

In [14]:

```

1  def sigmoid(Z):
2      return 1/(1+np.exp(-Z))
3
4  def relu(Z):
5      return np.maximum(0,Z)
6
7  def tanh(Z):
8      return np.tanh(Z)
9
10 def sigmoid_backward(dA, Z):
11     sig = sigmoid(Z)
12     return dA * sig * (1 - sig)
13
14 def relu_backward(dA, Z):
15     dZ = np.array(dA, copy = True) # pour ne pas effacer dA
16     dZ[Z <= 0] = 0;
17     return dZ;
18
19 def tanh_backward (dA, Z):
20     return 1- dA**2
21
22

```

La fonction suivante réalise la forward propagation mais uniquement d'un layer. Elle effectue donc le produit matriciel avec ajout du biais pour obtenir  $Z$ . Elle applique ensuite la fonction d'activation du layer.

In [15]:

```
1  def one_layer_forward_propagation(A_prev, W_curr, b_curr, activation="relu"):
2
3      # regression linéaire sur les entrées du layer
4      Z_curr = np.dot(W_curr, A_prev) + b_curr
5
6      # selection de la fonction d'activation à utiliser dans la couche
7      if activation == "relu":
8          activation_func = relu
9      elif activation == "sigmoid":
10         activation_func = sigmoid
11      elif activation == "tanh":
12         activation_func = tanh
13
14         # A est la sortie de la fonction d'activation
15         A=activation_func(Z_curr)
16         # retourne la fonction d'activation calculée et la matrice intermédiaire
17         return A, Z_curr
```

La fonction suivante fait la forward propagation sur tout le réseau. Elle sauvegarde aussi les valeurs de  $A$  et  $Z$  dans un dictionnaire python indexé par ces lettres afin de pouvoir les retrouver facilement lors de l'étape de backpropagation.

In [16]:

```

1  def forward_propagation(X, parameters, layers):
2
3      # Création d'un cache temporaire qui contient les valeurs intermédiaires
4      # utiles lors de la phase de backward. Le fait de les sauvegarder permet
5      # de ne pas les recalculer lors du backward
6      cache = {}
7
8      # A_curr correspond à la sortie du layer 0, i.e. les variables prédictive
9      A_curr = X
10
11     # iteration pour l'ensemble des couches du réseau
12     for idx, layer in enumerate(layers):
13         # La numérotation des couches commence à 1 (0 pour la couche des donn
14         layer_idx = idx + 1
15
16         # Récupération de l'activation de l'itération précédente
17         A_prev = A_curr
18
19         # Récupération du nom de la fonction d'activation du layer courant
20         activ_function_curr = layer["activation"]
21
22         # Récupération du W et du b du layer courant
23         W_curr = parameters["W" + str(layer_idx)]
24         b_curr = parameters["b" + str(layer_idx)]
25
26
27         # calcul de la fonction d'activation pour le layer courant
28         A_curr, Z_curr = one_layer_forward_propagation(A_prev, W_curr, b_curr
29
30         # Sauvegarde dans le cache pour la phase de backward
31         cache["A" + str(idx)] = A_prev
32         cache["Z" + str(layer_idx)] = Z_curr
33
34         # retourne le vecteur de prédiction à la sortie du réseau
35         # et un dictionnaire contenant toutes les valeurs intermédiaire
36         # pour faciliter la descente de gradient
37
38     return A_curr, cache

```

Le calcul de la fonction de coût (ici la cross entropy).  $\hat{y}$  correspond à la sortie du réseau après application de la forward propagation.  $y$  est la valeur réelle.

In [17]:

```

1  def cost_function(y_hat, y):
2      # calcul du coût (cross-entropy)
3      cost = (-y * np.log(y_hat) - (1 - y) * np.log(1 - y_hat)).mean()
4      return cost

```

Les deux fonctions suivantes permettent de calculer l'accuracy du modèle. Elles suivent le même principe que celles vues dans le notebook descente de gradient.

In [18]:

```

1  def convert_prob_into_class(probs):
2      probs = np.copy(probs) #pour ne pas perdre probs
3      probs[probs > 0.5] = 1
4      probs[probs <= 0.5] = 0
5      return probs
6
7  def accuracy(y_hat, y):
8      y_hat_ = convert_prob_into_class(y_hat)
9      return (y_hat_ == y).all(axis=0).mean()

```

Propagation d'un niveau. Tout d'abord nous considérons la descente de gradient sur un niveau. Elle applique tout d'abord la dérivée de la fonction d'activation passée en paramètre, dW, db et donc la dérivée de Z pour un layer L.

In [19]:

```

1  def one_layer_backward_propagation(dA_curr, W_curr, b_curr, Z_curr, A_prev, a
2
3      # nombres d'exemples venant de la fonction d'activation précédente
4      m = A_prev.shape[1]
5
6      # selection de la fonction d'activation à appliquer
7  if activation == "relu":
8      backward_activation_func = relu_backward
9  elif activation == "sigmoid":
10     backward_activation_func = sigmoid_backward
11  elif activation == "tanh":
12     backward_activation_func = tanh_backward
13
14     # calcul de la dérivée de la fonction d'activation
15     dZ_curr = backward_activation_func(dA_curr, Z_curr)
16
17     # dérivée de la matrice W
18     dW_curr = np.dot(dZ_curr, A_prev.T) / m
19     # dérivée du vecteur b
20     db_curr = np.sum(dZ_curr, axis=1, keepdims=True) / m
21     # dérivée de la matrice A_Prev
22     dA_prev = np.dot(W_curr.T, dZ_curr)
23
24     return dA_prev, dW_curr, db_curr

```

La propagation sur tout le réseau commence par calculer la dérivée de la fonction de coût :

$$\frac{\partial L}{\partial A} = \left( \frac{-y}{A} + \frac{(1-y)}{(1-A)} \right)$$

ou

$$- \left( \frac{y}{A} - \frac{(1-y)}{(1-A)} \right)$$

et boucle sur les autres niveaux. Le cache est utilisé pour récupérer les valeurs de A et de Z en fonction de leur layer et ce cache est utilisé pour sauvegarder dW et db qui seront utilisés lors de la descente de gradient.



In [20]:

```

1  def backward_propagation(y_hat, y, cache, parameters, layers):
2
3
4      # Création d'un cache temporaire qui contient les dérivées (gradients) po
5      # les différentes couches. Il est utilisé pour mettre à jour les paramètr
6      derivatives = {}
7
8      # nombre d'exemples
9      m = y.shape[1]
10
11     # pour garantir que y a la même forme que y_hat
12     y = y.reshape(y_hat.shape)
13
14     # Initialisation du calcul de la dérivée de la fonction de coût
15     # par rapport à A pour la couche L
16     dA_prev = - (np.divide(y, y_hat) - np.divide(1 - y, 1 - y_hat))
17
18     # parcours du réseau de la fonction finale vers celle d'entrée
19     for layer_idx_prev, layer in reversed(list(enumerate(layers))):
20         # Comme précédemment la numérotation des couches commence
21         # à 1. On ne modifie donc pas la couche 0
22         layer_idx_curr = layer_idx_prev + 1
23
24         # Récupération du nom de la fonction d'activation du layer courant
25         activ_function_curr = layer["activation"]
26
27         dA_curr = dA_prev
28
29         # Récupération dans le cache de la sortie précédente (A_prev)
30         # et de la matrice Z correspondant à l'application de la
31         # regression linéaire du niveau courant. Ceci évite de les recalculer
32         A_prev = cache["A" + str(layer_idx_prev)]
33         Z_curr = cache["Z" + str(layer_idx_curr)]
34
35         # Récupération dans parameters des valeurs de paramètres W et b du la
36         W_curr = parameters["W" + str(layer_idx_curr)]
37         b_curr = parameters["b" + str(layer_idx_curr)]
38
39         #application de la backwart propagation pour le layer afin d'avoir
40         #la valeur des dérivées (gradients)
41         dA_prev, dW_curr, db_curr = one_layer_backward_propagation(
42             dA_curr, W_curr, b_curr, Z_curr, A_prev, activ_function_curr)
43
44         # sauvegarde pour mettre à jour les paramètres
45         derivatives["dW" + str(layer_idx_curr)] = dW_curr
46         derivatives["db" + str(layer_idx_curr)] = db_curr
47
48     return derivatives

```

Application de la descente de gradient pour mettre à jour les paramètres. Comme tous les gradients ont été calculés précédemment (et sauvegardés dans un dictionnaire), il suffit de les appliquer. Ici la descente est faite à la manière d'une descente par lot (*batch gradient descent*) et peut être facilement modifiée en mini-batch gradient descent (C.f. notebook descente de gradient) avec optimisation.

In [21]:

```

1  ▼ def update(parameters, derivatives, layers, eta):
2
3      # Mise à jour des paramètres sur les différentes couches
4  ▼  for layer_idx, layer in enumerate(layers, 1):
5          parameters["w" + str(layer_idx)] -= eta * derivatives["dw" + str(layer_idx)]
6          parameters["b" + str(layer_idx)] -= eta * derivatives["db" + str(layer_idx)]
7
8      return parameters

```

La fonction train permet de lancer les différentes phases en fonction du nombre d'epochs. Elle retourne l'historique du coût et de l'accuracy pour pouvoir les afficher.

In [22]:

```

1  ▼ def fit(X, y, layers, epochs, eta):
2
3      # Initialisation des paramètres du réseau de neurones
4      parameters = init_layers(layers)
5
6      # sauvegarde historique coût et accuracy pour affichage
7      cost_history = []
8      accuracy_history = []
9
10     # Descente de gradient
11  ▼  for i in range(epochs):
12
13         # forward propagation
14         y_hat, cache = forward_propagation(X, parameters, layers)
15
16         # backward propagation - calcul des gradients
17         derivatives = backward_propagation(y_hat, y, cache, parameters, layers)
18
19         # Mise à jour des paramètres
20         parameters = update(parameters, derivatives, layers, eta)
21
22         # sauvegarde des historiques
23         current_cost = cost_function(y_hat, y)
24         cost_history.append(current_cost)
25         current_accuracy = accuracy(y_hat, y)
26         accuracy_history.append(current_accuracy)
27
28  ▼  if(i % 100 == 0):
29         print("Epoch : %s - cost : %.3f - accuracy : %.3f"%(i, float(current_cost), float(current_accuracy)))
30
31  return parameters, cost_history, accuracy_history

```

Premier test avec la configuration de l'exemple

In [23]:

```

1  def plot_histories (eta,epochs,cost_history,accuracy_history):
2      fig,ax = plt.subplots(figsize=(5,5))
3      ax.set_ylabel(r'$J(\theta)$')
4      ax.set_xlabel('Epochs')
5      ax.set_title(r"$\eta$ : {}".format(eta))
6      line1, = ax.plot(range(epochs),cost_history,label='Cost')
7      line2, = ax.plot(range(epochs),accuracy_history,label='Accuracy')
8      plt.legend(handler_map={line1: HandlerLine2D(numpoints=4)})

```

In [24]:

```

1  X,y = make_moons(n_samples=1000, noise=0.1)
2
3
4  validation_size=0.6 #40% du jeu de données pour le test
5
6  testsize= 1-validation_size
7  seed=30
8  # séparation jeu d'apprentissage et jeu de test
9  X_train,X_test,y_train,y_test=train_test_split(X,
10                                             y,
11                                             train_size=validation_size,
12                                             random_state=seed,
13                                             test_size=testsize)
14
15  # sauvegarde pour comparaison avec Keras
16  X_train_init=X_train
17  X_test_init=X_test
18  y_train_init=y_train
19  y_test_init=y_test
20
21  #transformation des données pour être au bon format
22  X_train=np.transpose(X_train)
23  X_test=np.transpose(X_test)
24  y_train=np.transpose(y_train.reshape((y_train.shape[0], 1)))
25  y_test=np.transpose(y_test.reshape((y_test.shape[0], 1)))

```

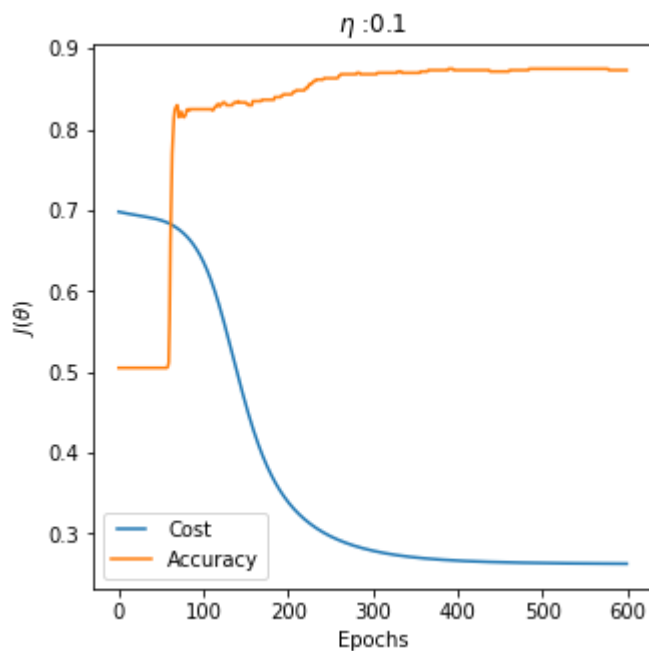
In [25]:

```

1  epochs = 600
2  eta = 0.1
3
4  ▼ parameters, cost_history, accuracy_history = fit(X_train,
5                                                    y_train,
6                                                    layers, epochs, eta)
7
8
9  # Calcul de l'accuracy sur le jeu de test
10 y_test_hat, _ = forward_propagation(X_test, parameters, layers)
11 accuracy_test = accuracy(y_test_hat, y_test)
12 print("Accuracy : %.3f"%accuracy_test)
13
14 plot_histories (eta, epochs, cost_history, accuracy_history)

```

Epoch : #0 - cost : 0.698 - accuracy : 0.505  
 Epoch : #100 - cost : 0.636 - accuracy : 0.825  
 Epoch : #200 - cost : 0.339 - accuracy : 0.843  
 Epoch : #300 - cost : 0.279 - accuracy : 0.868  
 Epoch : #400 - cost : 0.267 - accuracy : 0.873  
 Epoch : #500 - cost : 0.264 - accuracy : 0.875  
 Accuracy : 0.887



Test en modifiant le réseau pour voir si le résultat est meilleur. Attention au surapprentissage dans ce cas.

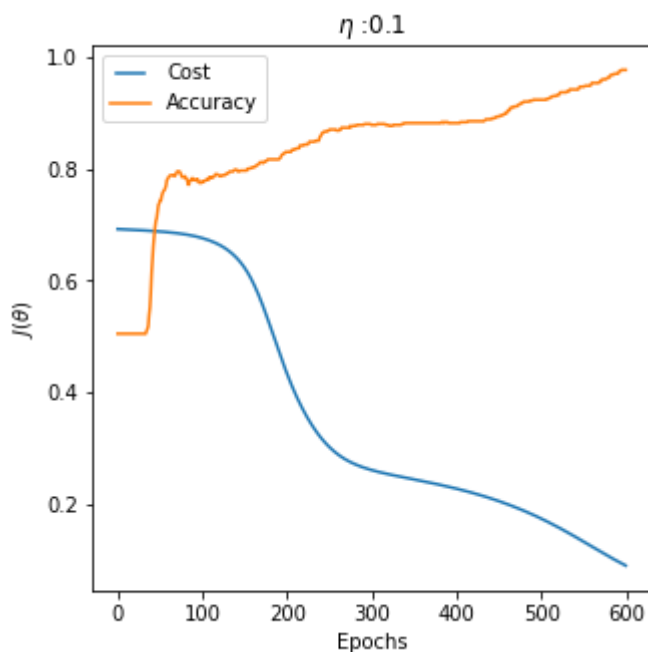
In [26]:

```

1  layers = [
2      {"input_dim": 2, "output_dim": 25, "activation": "relu"},
3      {"input_dim": 25, "output_dim": 25, "activation": "relu"},
4      {"input_dim": 25, "output_dim": 25, "activation": "relu"},
5      {"input_dim": 25, "output_dim": 1, "activation": "sigmoid"},
6  ]
7
8
9  epochs = 600
10 eta = 0.1
11 parameters, cost_history, accuracy_history = fit(X_train,
12                                                  y_train,
13                                                  layers, epochs, eta)
14
15 # Calcul de l'accuracy sur le jeu de test
16 y_test_hat, _ = forward_propagation(X_test, parameters, layers)
17 accuracy_test = accuracy(y_test_hat, y_test)
18 print("Accuracy : %.3f"%accuracy_test)
19
20 plot_histories (eta, epochs, cost_history, accuracy_history)

```

Epoch : #0 - cost : 0.692 - accuracy : 0.505  
Epoch : #100 - cost : 0.676 - accuracy : 0.778  
Epoch : #200 - cost : 0.432 - accuracy : 0.830  
Epoch : #300 - cost : 0.261 - accuracy : 0.878  
Epoch : #400 - cost : 0.228 - accuracy : 0.882  
Epoch : #500 - cost : 0.174 - accuracy : 0.923  
Accuracy : 0.983



Essai avec Keras pour voir les résultats.

In [27]:

```
1  ▼  #Construction du modèle
2  model = Sequential()
3  model.add(Dense(25, input_dim=2, activation='relu'))
4  model.add(Dense(25, activation='relu'))
5  model.add(Dense(25, activation='relu'))
6  model.add(Dense(1, activation='sigmoid'))
7  epochs = 600
8  eta = 0.1
9
10 gd = SGD(lr=eta)
11 model.compile(loss='binary_crossentropy', optimizer="sgd", metrics=['accuracy'])
12
13
14 # Training
15 history = model.fit(X_train_init, y_train_init, epochs=600, verbose=0)
16
17
```

In [28]:

```
1  y_test_hat = model.predict_classes(X_test_init)
2  accuracy_test = accuracy_score(y_test_init, y_test_hat)
3  print("Accuracy pour Keras : %.3f"%accuracy_test)
```

Accuracy pour Keras : 1.000

## Implémentation avec une classe simple

Dans un premier temps, nous créons une classe simple qui reprend les fonctions précédentes et dans laquelle nous ajoutons simplement les fonctions d'activation Relu, LeakyRelu, Tanh, Sigmoid et Softmax. Cette classe permet de faire de la classification multi-classes.

Récupération des librairies utiles

In [29]:

```
1  # importations utiles
2  import numpy as np
3  from sklearn.datasets import make_moons
4  from sklearn.datasets import make_circles
5  from sklearn import datasets
6  from sklearn.preprocessing import LabelEncoder
7  from sklearn.preprocessing import StandardScaler
8  from sklearn.preprocessing import OneHotEncoder
9  from matplotlib.colors import ListedColormap
10 import matplotlib.pyplot as plt
11 import matplotlib
12 from sklearn.model_selection import train_test_split
13 from sklearn import linear_model
14 from matplotlib.legend_handler import HandlerLine2D
15 import keras
16 from keras.models import Sequential
17 from keras.layers import Dense
18 from keras.utils import np_utils
19 from keras import regularizers
20 from keras.optimizers import SGD
21 from sklearn.metrics import accuracy_score
22 import pandas as pd
23 import seaborn as sns
24 import numpy as np
25 matplotlib.rcParams['figure.figsize'] = (6.0, 6.0) # pour avoir des figures d
```

Fonctions d'affichage.

In [30]:

```

1  def plot_histories (eta,epochs,cost_history,accuracy_history):
2      fig,ax = plt.subplots(figsize=(5,5))
3      ax.set_ylabel(r'$J(\theta)$')
4      ax.set_xlabel('Epochs')
5      ax.set_title(r"$\eta$ : {}".format(eta))
6      line1, = ax.plot(range(epochs),cost_history,label='Cost')
7      line2, = ax.plot(range(epochs),accuracy_history,label='Accuracy')
8      plt.legend(handler_map={line1: HandlerLine2D(numpoints=4)})
9
10 def plot_decision_boundary(func, X, y):
11     amin, bmin = X.min(axis=0) - 0.1
12     amax, bmax = X.max(axis=0) + 0.1
13     hticks = np.linspace(amin, amax, 101)
14     vticks = np.linspace(bmin, bmax, 101)
15
16     aa, bb = np.meshgrid(hticks, vticks)
17     ab = np.c_[aa.ravel(), bb.ravel()]
18     c = func(ab)
19     cc = c.reshape(aa.shape)
20
21     cm = plt.cm.RdBu
22     cm_bright = ListedColormap(['#FF0000', '#0000FF'])
23
24     fig, ax = plt.subplots()
25     contour = plt.contourf(aa, bb, cc, cmap=cm, alpha=0.8)
26
27     ax_c = fig.colorbar(contour)
28     ax_c.set_label("$P(y = 1)$")
29     ax_c.set_ticks([0, 0.25, 0.5, 0.75, 1])
30
31     plt.scatter(X[:, 0], X[:, 1], c=y, cmap=cm_bright)
32     plt.xlim(amin, amax)
33     plt.ylim(bmin, bmax)
34     plt.title("Decision Boundary")

```

Fonctions d'activations et dérivées des fonctions : sigmoid, tanh, relu, leakyrelu et softmax.



In [31]:

```

1  ▾ #fonctions utiles
2  ▾ def sigmoid(x):
3      return 1/(1 + np.exp(-x))
4
5  ▾ def sigmoid_prime(x):
6      return sigmoid(x)*(1.0 - sigmoid(x))
7
8  ▾ def tanh(x):
9      return np.tanh(x)
10
11 ▾ def tanh_prime(x):
12     return 1 - x ** 2
13
14 ▾ def relu(x):
15     return np.maximum(0,x)
16
17 ▾ def relu_prime(x):
18     x[x<=0] = 0
19     x[x>0] = 1
20     return x
21
22 ▾ def leakyrelu(x):
23     return np.maximum(0.01,x)
24
25 ▾ def leakyrelu_prime(x):
26     x[x<=0] = 0.01
27     x[x>0] = 1
28     return x
29
30 ▾ def softmax(x):
31     expx = np.exp(x - np.max(x))
32     return expx / expx.sum(axis=0, keepdims=True)

```

Nous définissons, à présent, une classe Layer. Son objectif est de pouvoir être utilisée de la même manière que dans Keras. Un layer peut être ajouté à notre réseau. Ce dernier va contenir le nombre de neurones de la couche ainsi que la fonction d'activation de la couche.

La première couche doit spécifier le nombre de neurones de la couche d'entrée.

Par exemple :

- `Layer(10,input=4,activation="leakyrelu")`

aura pour effet de créer un layer contenant 10 neurones, dont la couche d'entrée contient 4 neurones et pour lequel la fonction d'activation est leakyrelu.

Pour les autres couches il suffit de spécifier le nombre de couches et la fonction d'activation :

- `Layer(10,activation="softmax")`

ajoutera une couche contenant 10 neurones et dont la fonction d'activation est softmax (ici il s'agit de la dernière couche du réseau).

Les différents initialisations sont créées dans le constructeur de la classe Layer. La méthode `initParams` sera appelée par une méthode de la classe suivante et aura pour effet d'initialiser les valeurs de W et de b avec des nombres aléatoires. W est initialisé en prenant en compte l'optimisation de [He et al 2015] (Source:

<https://arxiv.org/pdf/1502.01852.pdf>) (<https://arxiv.org/pdf/1502.01852.pdf>).

In [32]:

```

1  class Layer:
2      def __init__(self,output,*args,**kwargs):
3
4          self.output = output # Number of neurons at layer i (current layer)
5          self.input = kwargs.get("input",None) # Number of neurons at layer i-
6          self.activ_function_curr = kwargs.get("activation",None) # Activation
7          self.parameters = {}
8          self.derivatives = {}
9          self.activation_func=None
10         if self.activ_function_curr == "relu":
11             self.activation_func = relu
12             self.backward_activation_func = relu_prime
13         elif self.activ_function_curr == "sigmoid":
14             self.activation_func = sigmoid
15             self.backward_activation_func = sigmoid_prime
16         elif self.activ_function_curr == "tanh":
17             self.activation_func = tanh
18             self.backward_activation_func = tanh_prime
19         elif self.activ_function_curr == "leakyrelu":
20             self.activation_func = leakyrelu
21             self.backward_activation_func = leakyrelu_prime
22         elif self.activ_function_curr == "softmax":
23             self.activation_func = softmax
24             self.backward_activation_func = softmax
25
26         def initParams(self):
27             # initialisation du dictionnaire de données parameters contenant W, A
28             seed=30
29             np.random.seed(seed)
30             self.parameters['W']=np.random.randn(self.output,self.input)*np.sqrt(
31             self.parameters['b']=np.random.randn(self.output,1)*0.1
32
33         def setW(self,matW):
34             self.parameters['W']=np.copy(matW)
35
36         def setA(self,matA):
37             self.parameters['A']=np.copy(matA)
38
39         def setZ(self,matZ):
40             self.parameters['Z']=np.copy(matZ)
41
42         def setB(self,matB):
43             self.parameters['b']=np.copy(matB)
44
45         def setdW(self,matdW):
46             self.parameters['dW']=np.copy(matdW)
47
48         def setdA(self,matdA):
49             self.parameters['dA']=np.copy(matdA)
50
51         def setdZ(self,matdZ):
52             self.parameters['dZ']=np.copy(matdZ)
53
54         def setdB(self,matdB):
55             self.parameters['dB']=np.copy(matdB)
56

```

Définition de la classe MyNeuralNetwork.

L'ajout de couche se fait à l'aide de la méthode *addLayer*. Cette fonction permet de mettre à jour les objets de la classe *Layer*.

Par exemple :

```
network = MyNeuralNetwork()

network.addLayer(Layer(3,input=2,activation="relu"))
network.addLayer(Layer(1,activation="sigmoid"))
```

permet de créer un réseau contenant 2 couches. La première est composée de 3 neurones, la couche d'entrée est composée de 2 neurones, elle utilise la fonction d'activation relu. La seconde couche, dernière couche, contient 1 neurone et la fonction d'activation est une sigmoid.

La classe est assez similaire à l'approche par fonctions. La principale différence est dans les fonctions forward et backward qui traitent différemment le cas des fonctions d'activations. Ici la descente de gradient se fait par mini-batch : d'où appel à la fonction *next\_batch* vue dans le notebook descente de gradient qui retourne un batch de la taille de batchsize.

La classe contient différentes méthodes supplémentaires :

- *info* qui donne les informations sur les différentes couches (nombre de neurones par couche, valeurs des W et des b lors de l'initialisation)
- *set\_parametersW\_b (numlayer,matX,matb)* qui permet de mettre des valeurs d'initialisation aux W et b pour une couche. Les valeurs doivent être sous la forme d'une matrice pour W et d'un vecteur pour b.
- *plot\_W\_b\_epoch (epoch,parameter\_history)* qui permet de connaître les valeurs de W (visualisation) à une epoch donnée. Elle utilise l'historique d'exécution du réseau.

Les différentes valeurs de A,Z, W, b ainsi que quelques statistiques pour chaque valeur d'epoch sont sauvegardées dans un tableau de dictionnaire.

L'entrainement du classifieur se fait à l'aide de la fonction *fit* :

Les premières valeurs correspondent à X et y, ensuite les paramètres sont (sans ordre particulier) :

- "epochs". Par défaut 20.
- "verbose"=(True|False). Par défaut False, qui affiche le coût et l'accuracy pour chaque epoch.
- "eta". Par défaut 0.01.
- "batchsize". Par défaut 32.

Les valeurs de retour sont :

- la liste des W et b à la fin de l'apprentissage (retour de la variable *layer*).
- les historiques des cost et accuracy pour chaque epoch afin de faire de pouvoir afficher l'historique.
- l'historique des A, Z, b, W pour chaque epoch.

Exemple d'utilisation :

```
layers,cost_history,accuracy_history,parameter_history=network.fit(X_train, y_train, verbose=True, epochs=100, eta=0.1, batchsize=128)
```

permet de lancer le classifieur avec comme paramètres X\_train et y\_train, en affichant les valeurs de cost et d'accuracy pour chaque epoch (verbose=True), sur 100 epochs (epochs=100), avec un learning rate de 0.1 (eta=0.1) et avec des batchsize de 128 (batchsize=128).

In [33]:

```

1  class MyNeuralNetwork:
2      def __init__(self):
3          self.nbLayers=0
4          self.layers=[]
5
6      def info(self):
7          print("Content of the network:");
8          j=0;
9          for i in range(len(self.layers)):
10             print("Layer n° ",i," => ")
11             print("\tInput ", self.layers[i].input,
12                   "\tOutput", self.layers[i].output)
13             if (i != 0):
14                 print("\tActivation Function",self.layers[i].activation_func)
15                 print("\tW", self.layers[i].parameters['W'].shape,self.layers[i].parameters['b'].shape)
16                 print("\tb", self.layers[i].parameters['b'].shape,self.layers[i].parameters['b'].shape)
17
18
19      def addLayer(self,layer):
20          self.nbLayers += 1;
21          if (self.nbLayers==1):
22              # this is the first layer so adding a layer 0
23              layerZero=Layer(layer.input)
24              self.layers.append(layerZero)
25
26          self.layers.append(layer)
27          self.layers[self.nbLayers].input=self.layers[self.nbLayers-1].output
28          self.layers[self.nbLayers].output=self.layers[self.nbLayers].output
29          layer.initParams()
30
31
32
33      def set_parametersW_b (self,numlayer,matX,matb):
34          self.layers[numlayer].parameters['W']=np.copy(matX)
35          self.layers[numlayer].parameters['b']=np.copy(matb)
36
37
38      def forward_propagation(self, X):
39          #Init predictive variables for the input layer
40          self.layers[0].setA(X)
41
42          #Propagation for all the layers
43          for l in range(1, self.nbLayers + 1):
44              # Compute Z
45              self.layers[l].setZ(np.dot(self.layers[l].parameters['W'],
46                                         self.layers[l-1].parameters['A'])+self.layers[l].parameters['b'])
47              # Applying the activation function of the layer to Z
48              self.layers[l].setA(self.layers[l].activation_func(self.layers[l].Z))
49
50
51      def cost_function(self,y):
52          return -(y*np.log(self.layers[self.nbLayers].parameters['A'])+1e-8)
53
54      def backward_propagation(self,y):
55          #calcul de dZ dW et db pour le dernier layer
56          self.layers[self.nbLayers].derivatives['dZ']=self.layers[self.nbLayers].derivatives['dA']
57          self.layers[self.nbLayers].derivatives['dW']=np.dot(self.layers[self.nbLayers].derivatives['dZ'],
58                                                                np.transpose(self.layers[self.nbLayers].parameters['W']))
59          m=self.layers[self.nbLayers].parameters['A'].shape[1]#égal au nombre

```

```

60 self.layers[self.nbLayers].derivatives['db']=np.sum(self.layers[self.nbLayers].derivatives,
61                                                         axis=1, keepdims=True)
62
63 #calcul de dz dW db pour les autres layers
64 for l in range(self.nbLayers-1,0,-1) :
65     self.layers[l].derivatives['dz']=np.dot(np.transpose(self.layers[l+1].derivatives['dz'],
66                                                         self.layers[l+1].derivatives['dz'].shape[1],
67                                                         self.layers[l+1].derivatives['dz'].shape[0]),
68     self.layers[l].derivatives["dW"]=np.dot(self.layers[l].derivatives["dz"],
69     np.transpose(self.layers[l+1].parameters["W"], self.layers[l+1].parameters["W"].shape[1],
70     self.layers[l+1].parameters["W"].shape[0])
71     m=self.layers[l+1].parameters['A'].shape[1]#égal au nombre de colonnes de la layer suivante
72     self.layers[l].derivatives['db']=np.sum(self.layers[l].derivatives['dz'], axis=1, keepdims=True)
73
74
75 def update_parameters(self, eta) :
76     for l in range(1,self.nbLayers+1) :
77         self.layers[l].parameters['W']-=eta*self.layers[l].derivatives['dW']
78         self.layers[l].parameters['b']-=eta*self.layers[l].derivatives['db']
79
80 def convert_prob_into_class(self,probs):
81     probs = np.copy(probs)#pour ne pas perdre probs, i.e. y_hat
82     probs[probs > 0.5] = 1
83     probs[probs <= 0.5] = 0
84     return probs
85
86 def plot_W_b_epoch (self,epoch,parameter_history):
87     mat=[]
88     max_size_layer=0
89     for l in range(1, self.nbLayers+1):
90         value=parameter_history[epoch]['W'+str(l)]
91         if (parameter_history[epoch]['W'+str(l)].shape[1]>max_size_layer):
92             max_size_layer=parameter_history[epoch]['W'+str(l)].shape[1]
93         mat.append(value)
94     figure=plt.figure(figsize=((self.nbLayers+1)*3,int (max_size_layer/2)))
95     for nb_w in range (len(mat)):
96         plt.subplot(1, len(mat), nb_w+1)
97         plt.matshow(mat[nb_w],cmap = plt.cm.gist_rainbow,fignum=False)
98         plt.colorbar()
99     thelegend="Epoch "+str(epoch)
100     plt.title (thelegend)
101
102 def accuracy(self,y_hat, y):
103     if self.layers[self.nbLayers].activation_func==softmax:
104         # si la fonction est softmax, les valeurs sont sur différentes dimensions
105         # il faut utiliser argmax avec axis=0 pour avoir un vecteur qui indique la classe
106         # où est la valeur maximale à la fois pour y_hat et pour y
107         # comme cela il suffit de comparer les deux vecteurs qui indiquent la classe
108         # dans quelle ligne se trouve le max
109         y_hat_encoded=np.copy(y_hat)
110         y_hat_encoded = np.argmax(y_hat_encoded, axis=0)
111         y_encoded=np.copy(y)
112         y_encoded=np.argmax(y_encoded, axis=0)
113         return (y_hat_encoded == y_encoded).mean()
114     # la dernière fonction d'activation n'est pas softmax.
115     # par exemple sigmoid pour une classification binaire
116     # il suffit de convertir la probabilité du résultat en classe
117     y_hat_ = self.convert_prob_into_class(y_hat)
118     return (y_hat_ == y).all(axis=0).mean()
119
120 def predict(self, x):

```

```

121     self.forward_propagation(x)
122     return self.layers[self.nbLayers].parameters['A']
123
124     def next_batch(self, X, y, batchsize):
125         # pour avoir X de la forme : 2 colonnes, m lignes (exemples) et égale
126         # cela permet de trier les 2 tableaux avec un indices de permutation
127         X=np.transpose(X)
128         y=np.transpose(y)
129
130         m=len(y)
131         # permutation aléatoire de X et y pour faire des batchs avec des valeurs
132         indices = np.random.permutation(m)
133         X = X[indices]
134         y = y[indices]
135         for i in np.arange(0, X.shape[0], batchsize):
136             # creation des batchs de taille batchsize
137             yield (X[i:i + batchsize], y[i:i + batchsize])
138     def fit(self, X, y, *args, **kwargs):
139         epochs=kwargs.get("epochs",20)
140         verbose=kwargs.get("verbose",False)
141         eta =kwargs.get("eta",0.01)
142         batchsize=kwargs.get("batchsize",32)
143     #def fit(self, X, y, epochs, eta = 0.01,batchsize=64) :
144         # sauvegarde historique coût et accuracy pour affichage
145         cost_history = []
146         accuracy_history = []
147         parameter_history = []
148         for i in range(epochs):
149             i+=1
150             # sauvegarde des coûts et accuracy par mini-batch
151             cost_batch = []
152             accuracy_batch = []
153             # Descente de gradient par mini-batch
154             for (batchX, batchy) in self.next_batch(X, y, batchsize):
155                 # Extraction et traitement d'un batch à la fois
156
157                 # mise en place des données au bon format
158                 batchX=np.transpose(batchX)
159                 if self.layers[self.nbLayers].activation_func==softmax:
160                     # la classification n'est pas binaire, y a utilisé one-hot
161                     # le batchy doit donc être transposé et le résultat doit
162                     # être sous la forme d'une matrice de taille batchy.shape[0], batchy.shape[1]
163
164                     batchy=np.transpose(batchy.reshape((batchy.shape[0], batchy.shape[1])))
165                 else:
166                     # il s'agit d'une classification binaire donc shape[1] n'est pas
167                     # pas
168                     batchy=np.transpose(batchy.reshape((batchy.shape[0], 1)))
169                 #batchy=np.transpose(batchy.reshape((batchy.shape[0], 1)))
170                 self.forward_propagation(batchX)
171                 self.backward_propagation(batchy)
172                 self.update_parameters(eta)
173
174                 # sauvegarde pour affichage
175                 current_cost=self.cost_function(batchy)
176                 cost_batch.append(current_cost)
177                 y_hat = self.predict(batchX)
178                 current_accuracy = self.accuracy(y_hat, batchy)
179                 accuracy_batch.append(current_accuracy)
180
181             # SaveStats on W, B as well as values for A,Z, W, b

```

```

182     save_values = {}
183     save_values["epoch"] = i
184     for l in range(1, self.nLayers+1):
185         save_values["layer"+str(l)] = l
186         save_values["Wmean"+str(l)] = np.mean(self.layers[self.nLayers-l])
187         save_values["Wmax"+str(l)] = np.amax(self.layers[self.nLayers-l])
188         save_values["Wmin"+str(l)] = np.amin(self.layers[self.nLayers-l])
189         save_values["Wstd"+str(l)] = np.std(self.layers[self.nLayers-l])
190         save_values["bmean"+str(l)] = np.mean(self.layers[self.nLayers-l])
191         save_values["bmax"+str(l)] = np.amax(self.layers[self.nLayers-l])
192         save_values["bmin"+str(l)] = np.amin(self.layers[self.nLayers-l])
193         save_values["bstd"+str(l)] = np.std(self.layers[self.nLayers-l])
194         # be careful A,Z,W and b must be copied otherwise it is a reference
195         save_values["A"+str(l)] = np.copy(self.layers[self.nLayers-l].parameters["A"])
196         save_values["Z"+str(l)] = np.copy(self.layers[self.nLayers-l].parameters["Z"])
197         save_values["W"+str(l)] = np.copy(self.layers[self.nLayers-l].parameters["W"])
198         save_values["b"+str(l)] = np.copy(self.layers[self.nLayers-l].parameters["b"])
199
200     parameter_history.append(save_values)
201     # sauvegarde de la valeur moyenne des coûts et de l'accuracy du batch
202     current_cost = np.average(cost_batch)
203     cost_history.append(current_cost)
204     current_accuracy = np.average(accuracy_batch)
205     accuracy_history.append(current_accuracy)
206
207     if(verbose == True):
208         print("Epoch : #s/%s - %s/%s - cost : %.4f - accuracy : %.4f" % (i, self.nLayers, current_cost, current_accuracy))
209
210     return self.layers, cost_history, accuracy_history, parameter_history
211

```

Utilisation de la classe MyNeural Network pour une classification binaire.

In [34]:

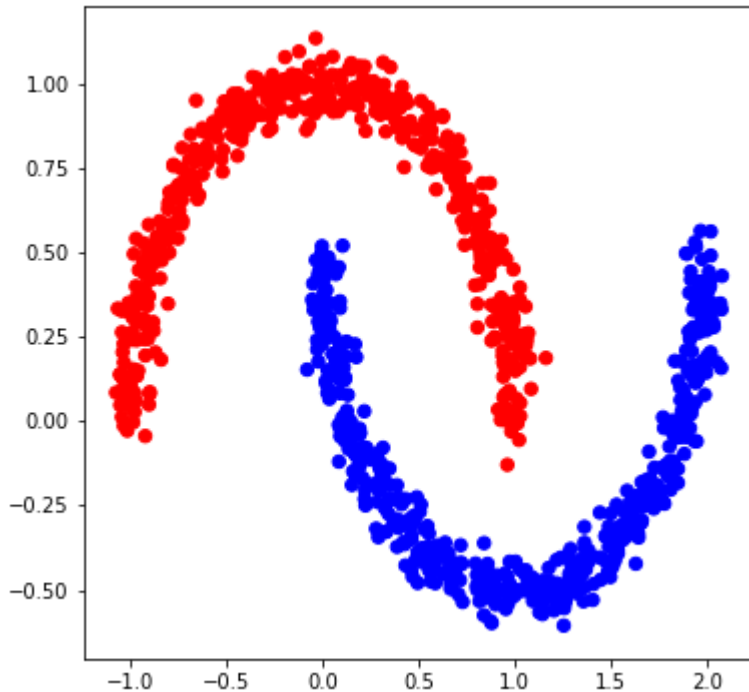
```

1 X, y = make_moons(n_samples=1000, noise=0.05, random_state=0)
2 cm_bright = ListedColormap(['#FF0000', '#0000FF'])
3 plt.scatter(X[:,0], X[:,1], s=40, c=y, cmap=cm_bright)#cmap=plt.cm.PiYG)

```

Out[34]:

&lt;matplotlib.collections.PathCollection at 0x12d1a3748&gt;



Création d'un jeu de données d'apprentissage et de test.

In [35]:

```

1  ▼ #création d'un jeu d'apprentissage et de test
2
3  validation_size=0.6 #40% du jeu de données pour le test
4
5  testsize= 1-validation_size
6  seed=30
7  # séparation jeu d'apprentissage et jeu de test
8  ▼ X_train,X_test,y_train,y_test=train_test_split(X,
9                                              y,
10                                             train_size=validation_size,
11                                             random_state=seed,
12                                             test_size=testsize)

```

Attention, les variables prédictives  $X$  doivent être sous la forme  $m$  colonnes où  $m$  est le nombre d'exemples d'apprentissage et  $l$  lignes où  $l$  est le nombre de variables prédictives. Dans notre exemple du notebook, il y a deux variables prédictives donc il faut passer que  $X$  soit de la forme  $m$  colonnes et 2 lignes.



In [36]:

```

1  #transformation des données pour être au bon format
2  # X_train est de la forme : 2 colonnes, m lignes (exemples)
3  # y_train est de la forme : m colonnes, 1 ligne
4
5  # La transposée de X_train est de la forme : m colonnes (exemples), 2 lignes
6  X_train=np.transpose(X_train)
7
8  # y_train est forcé pour être un tableau à 1 ligne contenant m colonnes
9  y_train=np.transpose(y_train.reshape((y_train.shape[0], 1)))
10
11 # mêmes traitements pour le jeu de test
12 X_test=np.transpose(X_test)
13 y_test=np.transpose(y_test.reshape((y_test.shape[0], 1)))
14

```

Création du réseau et affichage des informations.

In [37]:

```

1  network = MyNeuralNetwork()
2
3  network.addLayer(Layer(3,input=2,activation="relu"))
4  network.addLayer(Layer(1,activation="sigmoid"))
5
6  network.info()
7
8

```

Content of the network:

```

Layer n° 0 =>
    Input  None      Output 2
Layer n° 1 =>
    Input  2         Output 3
    Activation Function <function relu at 0x12e6941e0>
    W (3, 2) [[-1.26405266  1.52790535]
 [-0.97071094  0.47055962]
 [-0.10069672  0.30379318]]
    b (3, 1) [[-0.17259624]
 [ 0.15850954]
 [ 0.01342966]]
Layer n° 2 =>
    Input  3         Output 1
    Activation Function <function sigmoid at 0x12e6940d0>
    W (1, 3) [[-1.03209468  1.2475295  -0.79258216]]
    b (1, 1) [[0.04705596]]

```

Entraînement du classifieur et prédiction.

In [38]:

```

1  epochs = 100
2  eta = 0.01
3  batchsize=32
4
5  #Entraînement du classifieur
6  layers,cost_history,accuracy_history,parameter_history=network.fit(X_train, y
7
8
9  #Prédiction
10 y_pred=network.predict(X_test)
11 accuracy_test = network.accuracy(y_pred, y_test)
12 print("Accuracy test: %.3f"%accuracy_test)

```

```

Epoch : #1/100 - 600/600 - cost : 0.5145 - accuracy : 0.8947
Epoch : #2/100 - 600/600 - cost : 0.4712 - accuracy : 0.8931
Epoch : #3/100 - 600/600 - cost : 0.4599 - accuracy : 0.9019
Epoch : #4/100 - 600/600 - cost : 0.4503 - accuracy : 0.8920
Epoch : #5/100 - 600/600 - cost : 0.4434 - accuracy : 0.8893
Epoch : #6/100 - 600/600 - cost : 0.4366 - accuracy : 0.8958
Epoch : #7/100 - 600/600 - cost : 0.4290 - accuracy : 0.8936
Epoch : #8/100 - 600/600 - cost : 0.4215 - accuracy : 0.8920
Epoch : #9/100 - 600/600 - cost : 0.4162 - accuracy : 0.8920
Epoch : #10/100 - 600/600 - cost : 0.4098 - accuracy : 0.8969
Epoch : #11/100 - 600/600 - cost : 0.4052 - accuracy : 0.8914
Epoch : #12/100 - 600/600 - cost : 0.3981 - accuracy : 0.8947
Epoch : #13/100 - 600/600 - cost : 0.3920 - accuracy : 0.8964
Epoch : #14/100 - 600/600 - cost : 0.3772 - accuracy : 0.8931
Epoch : #15/100 - 600/600 - cost : 0.3164 - accuracy : 0.8871
Epoch : #16/100 - 600/600 - cost : 0.2832 - accuracy : 0.8832
Epoch : #17/100 - 600/600 - cost : 0.2665 - accuracy : 0.8854
Epoch : #18/100 - 600/600 - cost : 0.2568 - accuracy : 0.8876
Epoch : #19/100 - 600/600 - cost : 0.2531 - accuracy : 0.8876
Epoch : #20/100 - 600/600 - cost : 0.2449 - accuracy : 0.8882
Epoch : #21/100 - 600/600 - cost : 0.2415 - accuracy : 0.8849
Epoch : #22/100 - 600/600 - cost : 0.2371 - accuracy : 0.8871
Epoch : #23/100 - 600/600 - cost : 0.2348 - accuracy : 0.8904
Epoch : #24/100 - 600/600 - cost : 0.2325 - accuracy : 0.8914
Epoch : #25/100 - 600/600 - cost : 0.2291 - accuracy : 0.8887
Epoch : #26/100 - 600/600 - cost : 0.2279 - accuracy : 0.8898
Epoch : #27/100 - 600/600 - cost : 0.2260 - accuracy : 0.8936
Epoch : #28/100 - 600/600 - cost : 0.2227 - accuracy : 0.8865
Epoch : #29/100 - 600/600 - cost : 0.2226 - accuracy : 0.8860
Epoch : #30/100 - 600/600 - cost : 0.2218 - accuracy : 0.8882
Epoch : #31/100 - 600/600 - cost : 0.2204 - accuracy : 0.8904
Epoch : #32/100 - 600/600 - cost : 0.2201 - accuracy : 0.8920
Epoch : #33/100 - 600/600 - cost : 0.2157 - accuracy : 0.8898
Epoch : #34/100 - 600/600 - cost : 0.2129 - accuracy : 0.8947
Epoch : #35/100 - 600/600 - cost : 0.2133 - accuracy : 0.8914
Epoch : #36/100 - 600/600 - cost : 0.2113 - accuracy : 0.8871
Epoch : #37/100 - 600/600 - cost : 0.2111 - accuracy : 0.8953
Epoch : #38/100 - 600/600 - cost : 0.2117 - accuracy : 0.8871
Epoch : #39/100 - 600/600 - cost : 0.2106 - accuracy : 0.8909
Epoch : #40/100 - 600/600 - cost : 0.2113 - accuracy : 0.8788
Epoch : #41/100 - 600/600 - cost : 0.2106 - accuracy : 0.8860
Epoch : #42/100 - 600/600 - cost : 0.2121 - accuracy : 0.8920
Epoch : #43/100 - 600/600 - cost : 0.2098 - accuracy : 0.8942
Epoch : #44/100 - 600/600 - cost : 0.2059 - accuracy : 0.8942
Epoch : #45/100 - 600/600 - cost : 0.2095 - accuracy : 0.8865
Epoch : #46/100 - 600/600 - cost : 0.2097 - accuracy : 0.8893

```

```
Epoch : #47/100 - 600/600 - cost : 0.2124 - accuracy : 0.8920
Epoch : #48/100 - 600/600 - cost : 0.2077 - accuracy : 0.8904
Epoch : #49/100 - 600/600 - cost : 0.2054 - accuracy : 0.8849
Epoch : #50/100 - 600/600 - cost : 0.2077 - accuracy : 0.8931
Epoch : #51/100 - 600/600 - cost : 0.2029 - accuracy : 0.8909
Epoch : #52/100 - 600/600 - cost : 0.2055 - accuracy : 0.8893
Epoch : #53/100 - 600/600 - cost : 0.2077 - accuracy : 0.8882
Epoch : #54/100 - 600/600 - cost : 0.2041 - accuracy : 0.8887
Epoch : #55/100 - 600/600 - cost : 0.2049 - accuracy : 0.8914
Epoch : #56/100 - 600/600 - cost : 0.2054 - accuracy : 0.8854
Epoch : #57/100 - 600/600 - cost : 0.2035 - accuracy : 0.8953
Epoch : #58/100 - 600/600 - cost : 0.2025 - accuracy : 0.8975
Epoch : #59/100 - 600/600 - cost : 0.2023 - accuracy : 0.8909
Epoch : #60/100 - 600/600 - cost : 0.1946 - accuracy : 0.8958
Epoch : #61/100 - 600/600 - cost : 0.2010 - accuracy : 0.8991
Epoch : #62/100 - 600/600 - cost : 0.2036 - accuracy : 0.8942
Epoch : #63/100 - 600/600 - cost : 0.2021 - accuracy : 0.8991
Epoch : #64/100 - 600/600 - cost : 0.2015 - accuracy : 0.8920
Epoch : #65/100 - 600/600 - cost : 0.1993 - accuracy : 0.8969
Epoch : #66/100 - 600/600 - cost : 0.1980 - accuracy : 0.8958
Epoch : #67/100 - 600/600 - cost : 0.2004 - accuracy : 0.8986
Epoch : #68/100 - 600/600 - cost : 0.1924 - accuracy : 0.9013
Epoch : #69/100 - 600/600 - cost : 0.2045 - accuracy : 0.9030
Epoch : #70/100 - 600/600 - cost : 0.1974 - accuracy : 0.9008
Epoch : #71/100 - 600/600 - cost : 0.1946 - accuracy : 0.9013
Epoch : #72/100 - 600/600 - cost : 0.1959 - accuracy : 0.9068
Epoch : #73/100 - 600/600 - cost : 0.1928 - accuracy : 0.9046
Epoch : #74/100 - 600/600 - cost : 0.1950 - accuracy : 0.8991
Epoch : #75/100 - 600/600 - cost : 0.1926 - accuracy : 0.9030
Epoch : #76/100 - 600/600 - cost : 0.1940 - accuracy : 0.9117
Epoch : #77/100 - 600/600 - cost : 0.1913 - accuracy : 0.9062
Epoch : #78/100 - 600/600 - cost : 0.1933 - accuracy : 0.9123
Epoch : #79/100 - 600/600 - cost : 0.1914 - accuracy : 0.9019
Epoch : #80/100 - 600/600 - cost : 0.1876 - accuracy : 0.9112
Epoch : #81/100 - 600/600 - cost : 0.1894 - accuracy : 0.9183
Epoch : #82/100 - 600/600 - cost : 0.1868 - accuracy : 0.9150
Epoch : #83/100 - 600/600 - cost : 0.1888 - accuracy : 0.9189
Epoch : #84/100 - 600/600 - cost : 0.1895 - accuracy : 0.9101
Epoch : #85/100 - 600/600 - cost : 0.1866 - accuracy : 0.9117
Epoch : #86/100 - 600/600 - cost : 0.1873 - accuracy : 0.9254
Epoch : #87/100 - 600/600 - cost : 0.1893 - accuracy : 0.9183
Epoch : #88/100 - 600/600 - cost : 0.1831 - accuracy : 0.9139
Epoch : #89/100 - 600/600 - cost : 0.1841 - accuracy : 0.9183
Epoch : #90/100 - 600/600 - cost : 0.1847 - accuracy : 0.9249
Epoch : #91/100 - 600/600 - cost : 0.1891 - accuracy : 0.9221
Epoch : #92/100 - 600/600 - cost : 0.1914 - accuracy : 0.9260
Epoch : #93/100 - 600/600 - cost : 0.1834 - accuracy : 0.9205
Epoch : #94/100 - 600/600 - cost : 0.1821 - accuracy : 0.9232
Epoch : #95/100 - 600/600 - cost : 0.1823 - accuracy : 0.9243
Epoch : #96/100 - 600/600 - cost : 0.1830 - accuracy : 0.9211
Epoch : #97/100 - 600/600 - cost : 0.1806 - accuracy : 0.9304
Epoch : #98/100 - 600/600 - cost : 0.1896 - accuracy : 0.9227
Epoch : #99/100 - 600/600 - cost : 0.1770 - accuracy : 0.9238
Epoch : #100/100 - 600/600 - cost : 0.1826 - accuracy : 0.9271
Accuracy test: 0.945
```

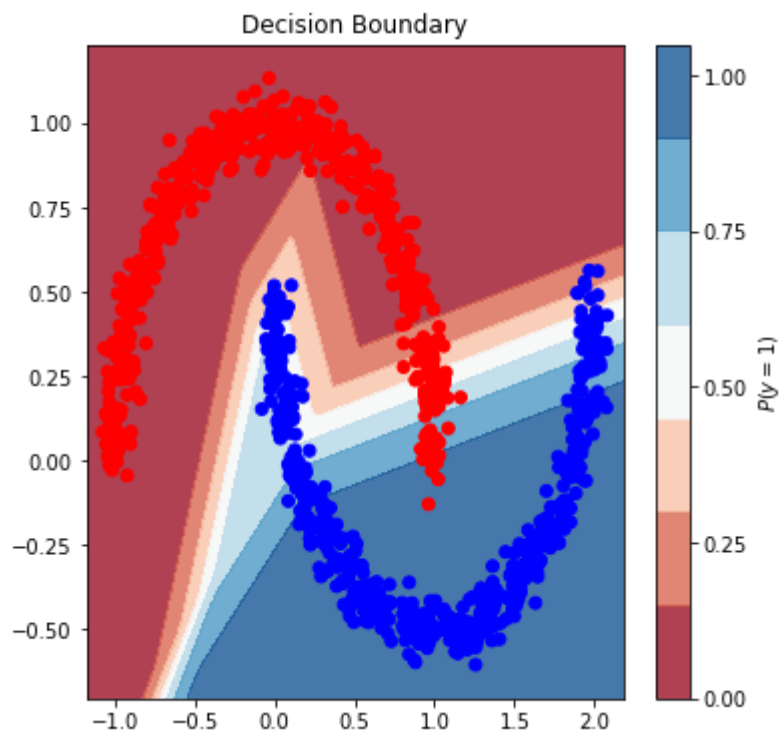
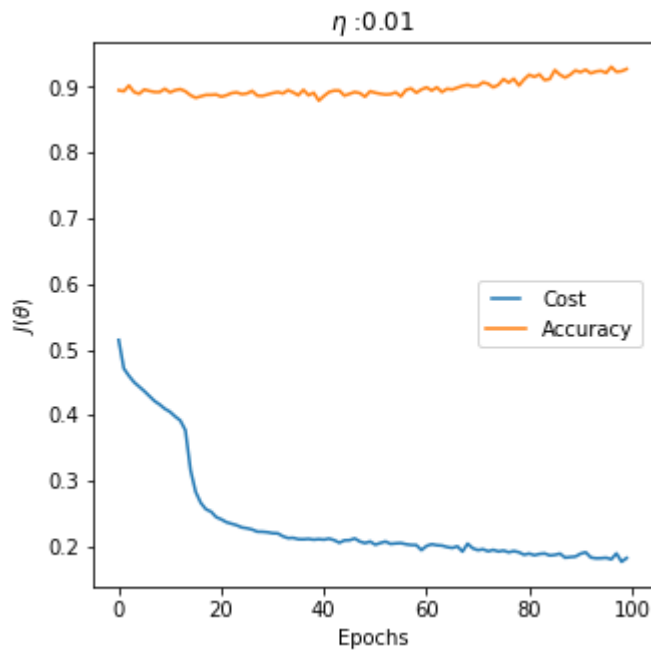
Affichage de l'historique de la fonction de coût et de l'accuracy ainsi que de la frontière de décision.

In [39]:

```

1  # Affichage des historiques
2  plot_histories (eta,epochs,cost_history,accuracy_history)
3  # Affichage de la frontière de décision
4  plot_decision_boundary(lambda x: network.predict(np.transpose(x)), X, y)
5  #plot_decision_boundary(lambda x: clf.predict(np.transpose(x)), X, y)

```



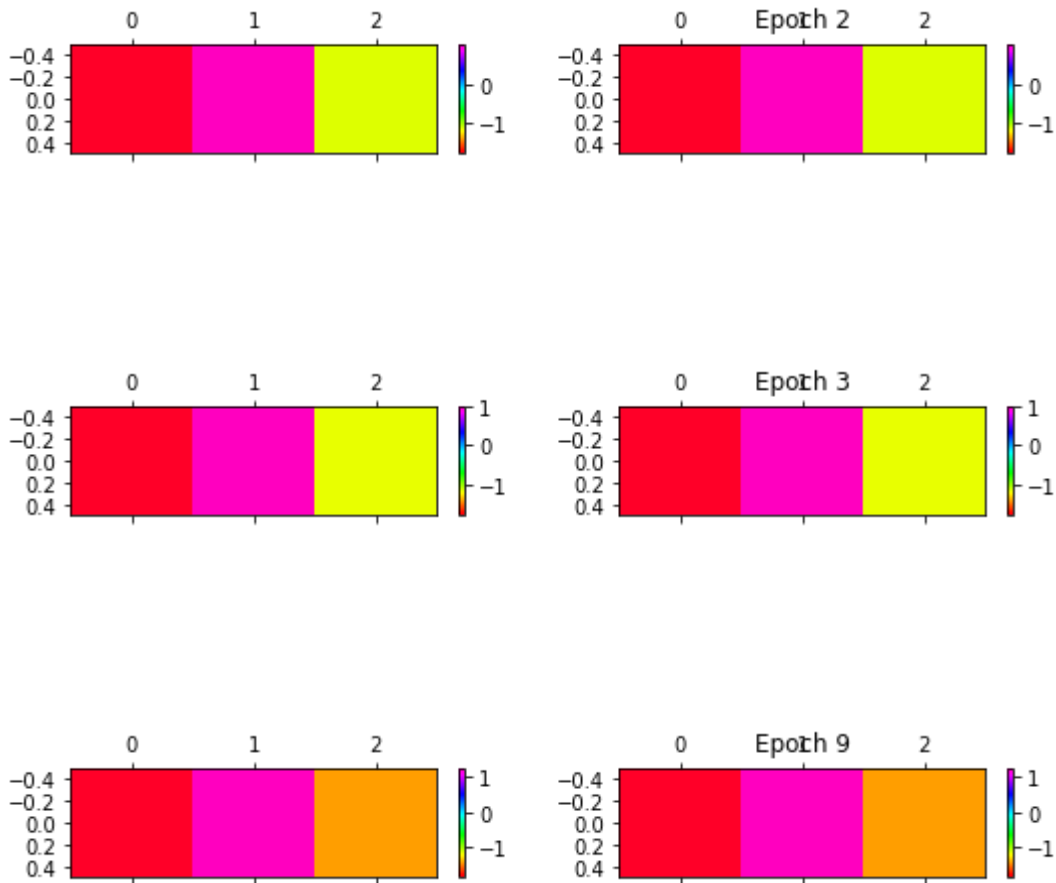
Affichage des valeurs de W pour les layers à différentes epochs.

In [40]:

```

1 layer_to_check=[2, 3, 9]
2 for i,e in enumerate(layer_to_check):
3     network.plot_W_b_epoch(e,parameter_history)
4

```



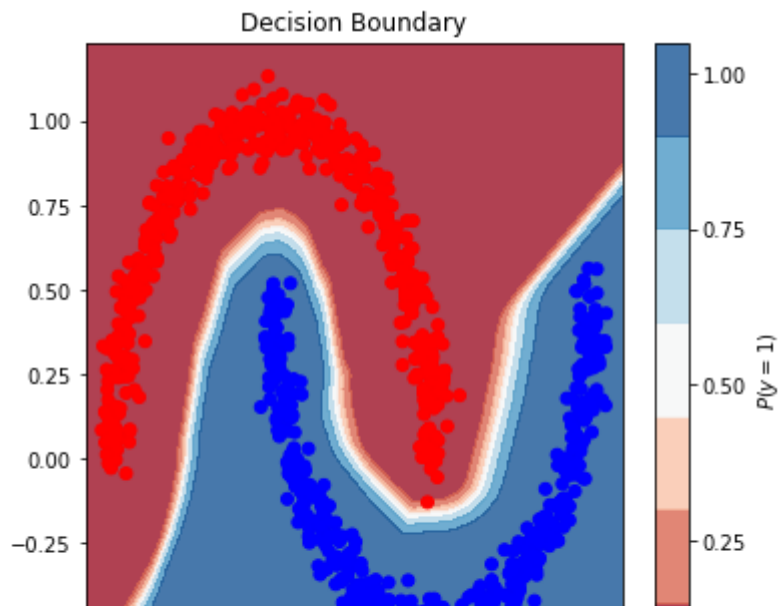
En changeant la configuration du réseau. Attention au surapprentissage.

In [41]:

```

1  network = MyNeuralNetwork()
2
3  network.addLayer(Layer(25,input=2,activation="relu"))
4  network.addLayer(Layer(25,activation="relu"))
5  network.addLayer(Layer(3,activation="relu"))
6  network.addLayer(Layer(1,activation="sigmoid"))
7
8  epochs = 100
9  eta = 0.01
10  batchsize=32
11
12  #Entraînement du classifieur
13  layers,cost_history,accuracy_history,parameter_history=network.fit(X_train, y
14
15
16  #Prédiction
17  y_pred=network.predict(X_test)
18  accuracy_test = network.accuracy(y_pred, y_test)
19  print("Accuracy test: %.3f"%accuracy_test)
20  # Affichage des historiques
21  plot_histories (eta,epochs,cost_history,accuracy_history)
22  # Affichage de la frontière de décision
23  plot_decision_boundary(lambda x: network.predict(np.transpose(x)), X, y)
24
25
26

```



Classification sur 3 classes à l'aide de softmax. Utilisation du jeu de données IRIS.

Sélection du jeu de données.

In [42]:

```

1 url="https://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data"
2 names = ['SepalLengthCm', 'SepalWidthCm',
3          'PetalLengthCm', 'PetalWidthCm',
4          'Species']
5
6 df = pd.read_csv(url, names=names)
7
8 # mélange des données
9 df=df.sample(frac=1).reset_index(drop=True)
10
11
12 array = df.values #nécessité de convertir le dataframe en numpy
13 #X matrice de variables prédictives - attention forcer le type à float
14 X = array[:,0:4].astype('float64')
15 #y vecteur de variable à prédire
16 y = array[:,4]
17

```

Première étape : normalisation des données

In [43]:

```

1 # normalisation de X
2 sc_X=StandardScaler()
3 X=sc_X.fit_transform(X)

```

Seconde étape : transformation des variables prédites à l'aide de OneHotEncoder : mettre 1 colonne par classe. Quand un exemple d'apprentissage correspond à la classe, il y a un 1 à la ligne et la colonne correspondante.

In [44]:

```

1 # Conversion de la variable à prédire via OneHotEncoder
2 # Dans IRIS il y a 3 classes -> création de 3 colonnes pour y
3 # 1 colonne correspond à 1 classe -> 1 si la ligne est du type de la classe
4 # 0 sinon
5
6 # Integer encode
7 label_encoder = LabelEncoder()
8 integer_encoded = label_encoder.fit_transform(y)
9
10 # binary encode
11 onehot_encoder = OneHotEncoder(sparse=False, categories='auto')
12 integer_encoded = integer_encoded.reshape(len(integer_encoded), 1)
13 y = onehot_encoder.fit_transform(integer_encoded)

```

Comme précédemment création du jeu d'apprentissage et de test. Attention, il faut mettre les variables prédictives et prédites au bon format.

In [45]:

```
1  ▼ # Jeu de test/apprentissage
2  validation_size=0.6 #40% du jeu de données pour le test
3
4  testsize= 1-validation_size
5  seed=30
6  # séparation jeu d'apprentissage et jeu de test
7  ▼ X_train,X_test,y_train,y_test=train_test_split(X,
8                                             y,
9                                             train_size=validation_size,
10                                          random_state=seed,
11                                          test_size=testsize)
12
13
14  #transformation des données pour être au bon format
15  # X_train est de la forme : n colonnes (variables à prédire après OneHotEncod
16  # y_train est de la forme : m colonnes, n lignes (variables à prédire après O
17
18  # La transposée de X_train est de la forme : m colonnes (exemples), n lignes
19  X_train=np.transpose(X_train)
20
21  # y_train est forcé pour être un tableau à 1 ligne contenant m colonnes
22  y_train=np.transpose(y_train.reshape((y_train.shape[0], y_train.shape[1])))
23
24  # mêmes traitements pour le jeu de test
25  X_test=np.transpose(X_test)
26  y_test=np.transpose(y_test.reshape((y_test.shape[0], y_test.shape[1])))
```

Création du réseau et lancement du classifieur.



In [46]:

```

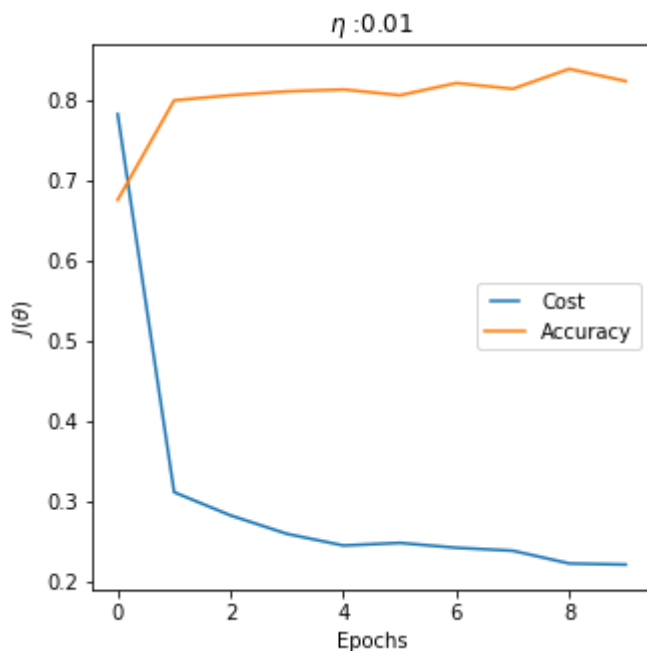
1  network = MyNeuralNetwork()
2
3  network.addLayer(Layer(10,input=4,activation="leakyrelu"))
4  network.addLayer(Layer(3,activation="softmax"))
5
6
7  epochs = 10
8  eta = 0.01
9  batchsize=128
10
11  #Entraînement du classifieur
12  layers,cost_history,accuracy_history,parameter_history=network.fit(X_train, y
13
14
15  #Prédiction
16  y_pred=network.predict(X_test)
17  accuracy_test = network.accuracy(y_pred, y_test)
18  print("Accuracy test: %.3f"%accuracy_test)
19
20  # Affichage des historiques
21  plot_histories (eta,epochs,cost_history,accuracy_history)

```

```

Epoch : #1/10 - 90/90 - cost : 0.7825 - accuracy : 0.6755
Epoch : #2/10 - 90/90 - cost : 0.3108 - accuracy : 0.7997
Epoch : #3/10 - 90/90 - cost : 0.2818 - accuracy : 0.8061
Epoch : #4/10 - 90/90 - cost : 0.2588 - accuracy : 0.8109
Epoch : #5/10 - 90/90 - cost : 0.2443 - accuracy : 0.8133
Epoch : #6/10 - 90/90 - cost : 0.2475 - accuracy : 0.8061
Epoch : #7/10 - 90/90 - cost : 0.2415 - accuracy : 0.8213
Epoch : #8/10 - 90/90 - cost : 0.2378 - accuracy : 0.8141
Epoch : #9/10 - 90/90 - cost : 0.2219 - accuracy : 0.8389
Epoch : #10/10 - 90/90 - cost : 0.2205 - accuracy : 0.8237
Accuracy test: 0.750

```



Essai sur le jeu de données MNIST. Un peu plus long car 60000 exemples d'apprentissage.

In [47]:

```

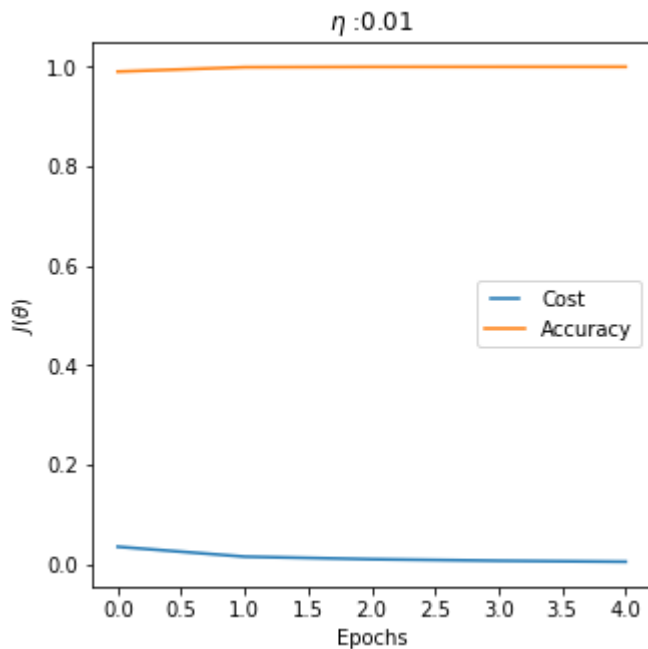
1  from tensorflow.keras.datasets import mnist
2  from keras.utils import to_categorical
3
4
5  (X_train_orig, y_train_orig), (X_test_orig, y_test_orig) = mnist.load_data()
6
7
8  y_tr_res = y_train_orig.reshape(60000, 1)
9  y_te_res = y_test_orig.reshape(10000, 1)
10 y_tr_T = to_categorical(y_tr_res, num_classes=10)
11 y_te_T = to_categorical(y_te_res, num_classes=10)
12 y_train = y_tr_T.T
13 #Y_train = Y_tr_res.T
14 y_test = y_te_T.T
15
16
17 X_train_flatten = X_train_orig.reshape(X_train_orig.shape[0], -1).T
18 #X_train_flatten=np.transpose(X_train)
19 X_test_flatten = X_test_orig.reshape(X_test_orig.shape[0], -1).T
20
21
22 X_train = X_train_flatten / 255.
23 X_test = X_test_flatten / 255.
24
25
26 network = MyNeuralNetwork()
27
28 network.addLayer(Layer(512,input=784,activation="relu"))
29 network.addLayer(Layer(10,activation="softmax"))
30
31 #network.info()
32 epochs = 5
33 eta = 0.01
34 batchsize=128
35 beta=0.9
36
37 #Entraînement du classifieur
38 layers,cost_history,accuracy_history,parameter_history=network.fit(X_train, y
39
40
41 #Prédiction
42 y_pred=network.predict(X_test)
43 accuracy_test = network.accuracy(y_pred, y_test)
44 print("Accuracy test: %.3f"%accuracy_test)
45
46 # Affichage des historiques
47 plot_histories (eta,epochs,cost_history,accuracy_history)
48

```

```

Epoch : #1/5 - 60000/60000 - cost : 0.0350 - accuracy : 0.9902
Epoch : #2/5 - 60000/60000 - cost : 0.0150 - accuracy : 0.9991
Epoch : #3/5 - 60000/60000 - cost : 0.0099 - accuracy : 0.9997
Epoch : #4/5 - 60000/60000 - cost : 0.0068 - accuracy : 0.9998
Epoch : #5/5 - 60000/60000 - cost : 0.0048 - accuracy : 0.9999
Accuracy test: 0.978

```



### Classe avec optimisation (momentum, adam)

Dans cette section, nous présentons une classe plus complète qui permet de faire des optimisations : momentum et adam.

Ces modifications sont répercutées sur la classe *Layer* et *MyNeuralNetwork*.

In [48]:

```

1  # importations utiles
2  import numpy as np
3  from sklearn.datasets import make_moons
4  from sklearn.datasets import make_circles
5  from sklearn import datasets
6  from sklearn.preprocessing import LabelEncoder
7  from sklearn.preprocessing import StandardScaler
8  from sklearn.preprocessing import OneHotEncoder
9  from matplotlib.colors import ListedColormap
10 import matplotlib.pyplot as plt
11 import matplotlib
12 from sklearn.model_selection import train_test_split
13 from sklearn import linear_model
14 from matplotlib.legend_handler import HandlerLine2D
15 import keras
16 from keras.models import Sequential
17 from keras.layers import Dense
18 from keras.utils import np_utils
19 from keras import regularizers
20 from keras.optimizers import SGD
21 from sklearn.metrics import accuracy_score
22 import pandas as pd
23 import seaborn as sns
24 import numpy as np
25 matplotlib.rcParams['figure.figsize'] = (6.0, 6.0) # pour avoir des figures d

```

Fonctions d'affichage inchangées.

In [49]:

```

1  def plot_histories (eta,epochs,cost_history,accuracy_history):
2      fig,ax = plt.subplots(figsize=(5,5))
3      ax.set_ylabel(r'$J(\theta)$')
4      ax.set_xlabel('Epochs')
5      ax.set_title(r"$\eta$ : {}".format(eta))
6      line1, = ax.plot(range(epochs),cost_history,label='Cost')
7      line2, = ax.plot(range(epochs),accuracy_history,label='Accuracy')
8      plt.legend(handler_map={line1: HandlerLine2D(numpoints=4)})
9
10 def plot_decision_boundary(func, X, y):
11     amin, bmin = X.min(axis=0) - 0.1
12     amax, bmax = X.max(axis=0) + 0.1
13     hticks = np.linspace(amin, amax, 101)
14     vticks = np.linspace(bmin, bmax, 101)
15
16     aa, bb = np.meshgrid(hticks, vticks)
17     ab = np.c_[aa.ravel(), bb.ravel()]
18     c = func(ab)
19     cc = c.reshape(aa.shape)
20
21     cm = plt.cm.RdBu
22     cm_bright = ListedColormap(['#FF0000', '#0000FF'])
23
24     fig, ax = plt.subplots()
25     contour = plt.contourf(aa, bb, cc, cmap=cm, alpha=0.8)
26
27     ax_c = fig.colorbar(contour)
28     ax_c.set_label("$P(y = 1)$")
29     ax_c.set_ticks([0, 0.25, 0.5, 0.75, 1])
30
31     plt.scatter(X[:, 0], X[:, 1], c=y, cmap=cm_bright)
32     plt.xlim(amin, amax)
33     plt.ylim(bmin, bmax)
34     plt.title("Decision Boundary")

```

Fonctions d'activations et dérivées inchangées.

In [50]:

```
1  ▼ #fonctions utiles
2  ▼ def sigmoid(x):
3      return 1/(1 + np.exp(-x))
4
5  ▼ def sigmoid_prime(x):
6      return sigmoid(x)*(1.0 - sigmoid(x))
7
8  ▼ def tanh(x):
9      return np.tanh(x)
10
11 ▼ def tanh_prime(x):
12     return 1 - x ** 2
13
14 ▼ def relu(x):
15     return np.maximum(0,x)
16
17 ▼ def relu_prime(x):
18     x[x<=0] = 0
19     x[x>0] = 1
20     return x
21
22 ▼ def leakyrelu(x):
23     return np.maximum(0.01,x)
24
25 ▼ def leakyrelu_prime(x):
26     x[x<=0] = 0.01
27     x[x>0] = 1
28     return x
29
30 ▼ def softmax(x):
31     expx = np.exp(x - np.max(x))
32     return expx / expx.sum(axis=0, keepdims=True)
33
34
```

Le constructeur et la méthode *initParams* ont été modifiés pour pouvoir prendre en compte les optimisations de descente de gradient.

In [51]:

```

1  class Layer:
2      def __init__(self,output,*args,**kwargs):
3          self.output = output # Nombre de neurones au layer i (actuel)
4          self.input = kwargs.get("input",None) # Nombre de neurones au layer i
5          self.activ_function_curr = kwargs.get("activation",None) # fonction d
6          self.parameters = {}
7          self.derivatives = {}
8          #for momentum
9          self.v = {}
10         #for adam
11         self.s = {}
12
13         self.activation_func=None
14         if self.activ_function_curr == "relu":
15             self.activation_func = relu
16             self.backward_activation_func = relu_prime
17         elif self.activ_function_curr == "sigmoid":
18             self.activation_func = sigmoid
19             self.backward_activation_func = sigmoid_prime
20         elif self.activ_function_curr == "tanh":
21             self.activation_func = tanh
22             self.backward_activation_func = tanh_prime
23         elif self.activ_function_curr == "leakyrelu":
24             self.activation_func = leakyrelu
25             self.backward_activation_func = leakyrelu_prime
26         elif self.activ_function_curr == "softmax":
27             self.activation_func = softmax
28             self.backward_activation_func = softmax
29
30         def initParams(self,optimizer):
31             seed=30
32             np.random.seed(seed)
33             self.parameters['W']=np.random.randn(self.output,self.input)*np.sqrt(
34             self.parameters['b']=np.random.randn(self.output,1)*0.1
35             if optimizer=="momentum":
36                 # Sauvegarde velocity, v, pour momentum
37                 self.v['dW'] = np.zeros_like(self.parameters['W'])
38                 self.v['db'] = np.zeros_like(self.parameters['b'])
39             elif optimizer=='adam':
40                 self.v['dW'] = np.zeros_like(self.parameters['W'])
41                 self.v['db'] = np.zeros_like(self.parameters['b'])
42                 self.s['dW'] = np.zeros_like(self.parameters['W'])
43                 self.s['db'] = np.zeros_like(self.parameters['b'])
44
45
46         def setW(self,matW):
47             #self.parameters['W']=matW
48             self.parameters['W']=np.copy(matW)
49
50         def setA(self,matA):
51             self.parameters['A']=np.copy(matA)
52
53         def setZ(self,matZ):
54             self.parameters['Z']=np.copy(matZ)
55
56         def setB(self,matB):
57             self.parameters['b']=np.copy(matB)
58
59         def setdW(self,matdW):

```

```

60         self.parameters[ 'dW' ]=np.copy(matdW)
61
62     def setdA(self,matdA):
63         self.parameters[ 'dA' ]=np.copy(matdA)
64
65     def setdZ(self,matdZ):
66         self.parameters[ 'dZ' ]=np.copy(matdZ)
67
68     def setdB(self,matdB):
69         self.parameters[ 'dB' ]=np.copy(matdB)
70

```

La structure de la classe *MyNeuralNetwork* est assez similaire à la précédente.

- Le constructeur prend par défaut *optimizer="bgd"* (*batch gradient descent*) (les valeurs possibles sont "*momentum*" ou "*adam*").
- La méthode *update* prend en compte le fait que l'optimisation puisse être de type *momentum* ou *adam*.

Par exemple :

```
network = MyNetwork(optimizer="adam")
```

aura pour effet de créer un réseau qui utilise l'optimisation *adam*.

Enfin lors de l'appel de la méthode *fit* les paramètres *beta*, *beta1* et *epsilon* doivent être spécifiées (en cas d'appel d'un *optimizer*. Ici les valeurs par défaut sont : *beta=0.9*, *beta2=0.999*, *epsilon=1e-8*.

Par exemple :

```
layers, cost_history, accuracy_history, parameter_history = network.fit(X_train, y_train, verbose=True, epochs=40,
eta=0.01, beta=0.8)
```

aura pour effet de lancer le classifieur sur *X\_train*, *y\_train*, en affichant le coût et l'accuracy (*verbose=True*), pour 40 epochs (*epochs=40*), avec un learning rate de 0.01 (*eta=0.01*) et un *beta* de 0.8 (*beta=0.8*).

In [52]:

```

1  class MyNeuralNetwork:
2      def __init__(self,optimizer="bgd"):
3          self.nbLayers=0
4          self.layers=[]
5          self.optimizer=optimizer
6          self.t=2
7
8      def info(self):
9          print("Content of the network:");
10         j=0;
11         for i in range(len(self.layers)):
12             print("Layer n° ",i," => ")
13             print("\tInput ", self.layers[i].input,
14                   "\tOutput", self.layers[i].output)
15             if (i != 0):
16                 print("\tActivation Function",self.layers[i].activation_func)
17                 print("\tW", self.layers[i].parameters['W'].shape,self.layers[i].parameters['b'].shape)
18                 print("\tb", self.layers[i].parameters['b'].shape,self.layers[i].parameters['W'].shape)
19
20
21         def addLayer(self,layer):
22             self.nbLayers += 1;
23             if (self.nbLayers==1):
24                 # this is the first layer so adding a layer 0
25                 layerZero=Layer(layer.input)
26                 self.layers.append(layerZero)
27
28             self.layers.append(layer)
29             self.layers[self.nbLayers].input=self.layers[self.nbLayers-1].output
30             self.layers[self.nbLayers].output=self.layers[self.nbLayers].output
31             layer.initParams(self.optimizer)
32
33
34
35         def set_parametersW_b (self,numlayer,matX,matb):
36             self.layers[numlayer].parameters['W']=np.copy(matX)
37             self.layers[numlayer].parameters['b']=np.copy(matb)
38
39
40         def forward_propagation(self, X):
41             #Initialisation des variables prédictives pour la couche d'entrée
42             self.layers[0].setA(X)
43
44             #Propagation pour tous les layers
45             for l in range(1, self.nbLayers + 1):
46                 # Calcul de Z
47                 self.layers[l].setZ(np.dot(self.layers[l].parameters['W'],
48                                             self.layers[l-1].parameters['A'])+self.layers[l].parameters['b'])
49                 # Application de la fonction d'activation à Z
50                 self.layers[l].setA(self.layers[l].activation_func(self.layers[l].Z))
51
52
53         def cost_function(self,y):
54             return -(y*np.log(self.layers[self.nbLayers].parameters['A'])+1e-8) -
55
56         def backward_propagation(self,y):
57             #calcul de dZ dW et db pour le dernier layer
58             self.layers[self.nbLayers].derivatives['dZ']=self.layers[self.nbLayers].output-y
59             self.layers[self.nbLayers].derivatives['dW']=np.dot(self.layers[self.nbLayers].derivatives['dZ'],self.layers[self.nbLayers].Z)

```



```

60                                     np.transpose(self
61 m=self.layers[self.nbLayers].parameters['A'].shape[1]#égal au nombre
62 self.layers[self.nbLayers].derivatives['db']=np.sum(self.layers[self
63                                     axis=1, keepdims=True)
64
65 #calcul de dZ dW db pour les autres layers
66 for l in range(self.nbLayers-1,0,-1) :
67     self.layers[l].derivatives['dZ']=np.dot(np.transpose(self.layers[
68                                     self.layers[l+1].derivatives['dZ
69
70     self.layers[l].derivatives["dW"]=np.dot(self.layers[l].derivative
71                                     np.transpose(self.layers[l-1].pa
72
73 m=self.layers[l-1].parameters['A'].shape[1]#égal au nombre de co
74 self.layers[l].derivatives['db']=np.sum(self.layers[l].derivative
75                                     axis=1, keepdims=True)
76
77
78 def update_parameters(self, eta,beta, beta2=0.999, epsilon=1e-8):
79     # Descente de gradient
80
81     if self.optimizer=="adam":
82         v_corrected = {} # Initialisation de la première estim
83         s_corrected = {} # Initialisation de la seconde estim
84
85     for l in range(1, self.nbLayers+1):
86         if self.optimizer=="momentum":
87             # Calcul de la vitesse
88             self.layers[l].v['dW'] = beta * self.layers[l].v['dW'] + (1 -
89             self.layers[l].v['db'] = beta * self.layers[l].v['db'] + (1 -
90
91             # Mise à jour des paramètres
92             self.layers[l].parameters['W'] -= eta*self.layers[l].v['dW']
93             self.layers[l].parameters['b'] -= eta*self.layers[l].v['db']
94         elif self.optimizer=="adam":
95             # Calcul de la vitesse
96             self.layers[l].v['dW'] = beta * self.layers[l].v['dW'] + (1 -
97             self.layers[l].v['db'] = beta * self.layers[l].v['db'] + (1 -
98
99             # Calcul de la première estimation du moment (correction du l
100 v_corrected['dW' + str(l)] = self.layers[l].v['dW'] / (1 - np
101 v_corrected['db' + str(l)] = self.layers[l].v['db'] / (1 - np
102
103             # déplacement moyen des gradients au carré
104             self.layers[l].s['dW'] = beta2 * self.layers[l].s['dW'] + (1
105             self.layers[l].s['db'] = beta2 * self.layers[l].s['db'] + (1
106
107             # Calcul de la seconde estimation du moment (correction du b
108 s_corrected["dW" + str(l)] = self.layers[l].s['dW'] / (1 - np
109 s_corrected["db" + str(l)] = self.layers[l].s['db'] / (1 - np
110
111             # Mise à jour des paramètres
112             self.layers[l].parameters['W'] -= eta*v_corrected['dW' + str
113             self.layers[l].parameters['b'] -= eta*v_corrected['db' + str
114
115         else: #descente par mini-lots
116             self.layers[l].parameters['W'] -= eta*self.layers[l].derivat
117             self.layers[l].parameters['b'] -= eta*self.layers[l].derivat
118
119
120 def convert_prob_into_class(self,probs):

```

```

121     probs = np.copy(probs) #pour ne pas perdre probs, i.e. y_hat
122     probs[probs > 0.5] = 1
123     probs[probs <= 0.5] = 0
124     return probs
125
126     def plot_W_b_epoch (self, epoch, parameter_history):
127         mat=[]
128         max_size_layer=0
129         for l in range(1, self.nbLayers+1):
130             value=parameter_history[epoch][ 'W'+str(l) ]
131             if (parameter_history[epoch][ 'W'+str(l) ].shape[1]>max_size_layer):
132                 max_size_layer=parameter_history[epoch][ 'W'+str(l) ].shape[1]
133             mat.append(value)
134         figure=plt.figure(figsize=((self.nbLayers+1)*3,int (max_size_layer/2)))
135         for nb_w in range (len(mat)):
136             plt.subplot(1, len(mat), nb_w+1)
137             plt.matshow(mat[nb_w], cmap = plt.cm.gist_rainbow, fignum=False)
138             plt.colorbar()
139         thelegend="Epoch "+str(epoch)
140         plt.title (thelegend)
141
142
143
144     def accuracy(self, y_hat, y):
145         if self.layers[self.nbLayers].activation_func==softmax:
146             # si la fonction est softmax, les valeurs sont sur différentes d
147             # il faut utiliser argmax avec axis=0 pour avoir un vecteur qui
148             # où est la valeur maximale à la fois pour y_hat et pour y
149             # comme cela il suffit de comparer les deux vecteurs qui indiquent
150             # dans quelle ligne se trouve le max
151             y_hat_encoded=np.copy(y_hat)
152             y_hat_encoded = np.argmax(y_hat_encoded, axis=0)
153             y_encoded=np.copy(y)
154             y_encoded=np.argmax(y_encoded, axis=0)
155             return (y_hat_encoded == y_encoded).mean()
156         # la dernière fonction d'activation n'est pas softmax.
157         # par exemple sigmoid pour une classification binaire
158         # il suffit de convertir la probabilité du résultat en classe
159         y_hat_ = self.convert_prob_into_class(y_hat)
160         return (y_hat_ == y).all(axis=0).mean()
161
162     def predict(self, x):
163         self.forward_propagation(x)
164         return self.layers[self.nbLayers].parameters['A']
165
166     def next_batch(self, X, y, batchsize):
167         # pour avoir X de la forme : 2 colonnes, m lignes (exemples) et égale
168         # cela permet de trier les 2 tableaux avec un indices de permutation
169         X=np.transpose(X)
170         y=np.transpose(y)
171
172         m=len(y)
173         # permutation aléatoire de X et y pour faire des batchs avec des valeurs
174         indices = np.random.permutation(m)
175         X = X[indices]
176         y = y[indices]
177         for i in np.arange(0, X.shape[0], batchsize):
178             # creation des batchs de taille batchsize
179             yield (X[i:i + batchsize], y[i:i + batchsize])
180
181

```

```

182 #def fit(self, X, y, epochs, verbose=True, eta = 0.01, batchsize=32, beta=0
183 def fit(self, X, y, *args, **kwargs):
184     epochs=kwargs.get("epochs", 20)
185     verbose=kwargs.get("verbose", False)
186     eta =kwargs.get("eta", 0.01)
187     batchsize=kwargs.get("batchsize", 32)
188     beta=kwargs.get("beta", 0.9)
189     beta2=kwargs.get("beta2", 0.999)
190     epsilon=kwargs.get("epsilon", 1e-8)
191     # sauvegarde historique coût, accuracy pour affichage
192     # parameter_history sauvegarde des statistiques, W et b pour toutes
193     # epochs afin de pouvoir étudier l'historique
194     cost_history = []
195     accuracy_history = []
196     parameter_history = []
197
198     for i in range(epochs):
199         i+=1
200         # sauvegarde des coûts et accuracy par mini-batch
201         cost_batch = []
202         accuracy_batch = []
203         # Gradient descent per mini-batch
204
205         for (batchX, batchy) in self.next_batch(X, y, batchsize):
206
207             # Extract and process one batch at a time
208
209             # Data must be at the appropriate format
210             batchX=np.transpose(batchX)
211
212             if self.layers[self.nbLayers].activation_func==softmax:
213                 # Not binary classification. one-hot-encoder has been used
214                 # the batchy must be transformed and the result
215                 # must be a matrice with size batchy.shape[1]
216
217                 batchy=np.transpose(batchy.reshape((batchy.shape[0], batchy.shape[1])))
218             else:
219                 # It is a binary classification so shape[1] does not exist
220                 batchy=np.transpose(batchy.reshape((batchy.shape[0], 1)))
221
222             self.forward_propagation(batchX)
223             self.backward_propagation(batchy)
224             if self.optimizer=="adam":
225                 self.t=self.t+1
226             self.update_parameters(eta, beta, beta2, epsilon)
227
228             # save for output
229             current_cost=self.cost_function(batchy)
230             cost_batch.append(current_cost)
231             y_hat = self.predict(batchX)
232             current_accuracy = self.accuracy(y_hat, batchy)
233             accuracy_batch.append(current_accuracy)
234
235             # SaveStats on W, B as well as values for A, Z, W, b
236             save_values = {}
237             save_values["epoch"]=i
238             for l in range(1, self.nbLayers+1):
239                 save_values["layer"+str(l)]=l
240                 save_values["Wmean"+ str(l)]=np.mean(self.layers[self.nbLayers+1-l].W)
241                 save_values["Wmax"+ str(l)]=np.amax(self.layers[self.nbLayers+1-l].W)

```

```

243         save_values["Wmin"+str(l)]=np.amin(self.layers[self.nbLayers-l].weights)
244         save_values["Wstd"+str(l)]=np.std(self.layers[self.nbLayers-l].weights)
245         save_values["bmean"+str(l)]=np.mean(self.layers[self.nbLayers-l].biases)
246         save_values["bmax"+str(l)]=np.amax(self.layers[self.nbLayers-l].biases)
247         save_values["bmin"+str(l)]=np.amin(self.layers[self.nbLayers-l].biases)
248         save_values["bstd"+str(l)]=np.std(self.layers[self.nbLayers-l].biases)
249         # be careful A,Z,W and b must be copied otherwise it is a reference
250         save_values["A"+str(l)]=np.copy(self.layers[self.nbLayers-l].activation)
251         save_values["Z"+str(l)]=np.copy(self.layers[self.nbLayers-l].z)
252         save_values["W"+str(l)]=np.copy(self.layers[self.nbLayers-l].weights)
253         save_values["b"+str(l)]=np.copy(self.layers[self.nbLayers-l].biases)
254
255         parameter_history.append(save_values)
256
257         # save avg value for cost and accuracy
258         current_cost=np.average(cost_batch)
259         cost_history.append(current_cost)
260         current_accuracy=np.average(accuracy_batch)
261         accuracy_history.append(current_accuracy)
262
263         if(verbose == True):
264             print("Epoch : #s/%s - %s/%s - cost : %.4f - accuracy : %.4f" % (l, nbEpochs,
265                                                                                   current_cost,
266                                                                                   current_accuracy))
267
268     return self.layers, cost_history, accuracy_history, parameter_history

```

Test sur un jeu de données

In [53]:

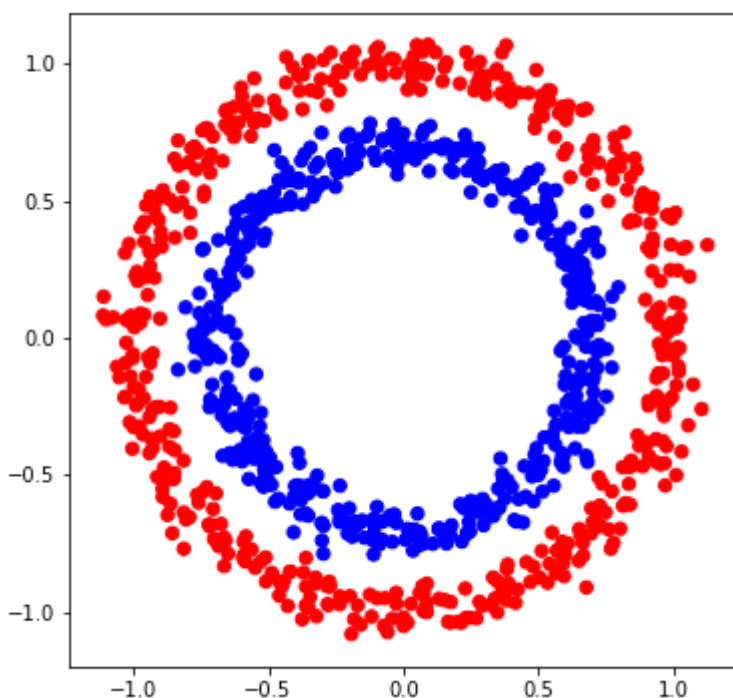
```

1 X, y = make_circles(n_samples=1000, noise=0.05, factor=0.7, random_state=0)
2 cm_bright = ListedColormap(['#FF0000', '#0000FF'])
3 plt.scatter(X[:,0], X[:,1], s=40, c=y, cmap=cm_bright)#cmap=plt.cm.PiYG)

```

Out[53]:

<matplotlib.collections.PathCollection at 0x12ce94fd0>



Préparation des données et du jeu d'apprentissage et de test.

In [54]:

```
1 validation_size=0.6 #40% du jeu de données pour le test
2
3 testsize= 1-validation_size
4 seed=30
5 # séparation jeu d'apprentissage et jeu de test
6 X_train,X_test,y_train,y_test=train_test_split(X,
7                                               y,
8                                               train_size=validation_size,
9                                               random_state=seed,
10                                              test_size=testsize)
11
12
13 #transformation des données pour être au bon format
14 # X_train est de la forme : 2 colonnes, m lignes (exemples)
15 # y_train est de la forme : m colonnes, 1 ligne
16
17 # La transposée de X_train est de la forme : m colonnes (exemples), 2 lignes
18 X_train=np.transpose(X_train)
19
20 # y_train est forcé pour être un tableau à 1 ligne contenant m colonnes
21 y_train=np.transpose(y_train.reshape((y_train.shape[0], 1)))
22
23 # mêmes traitements pour le jeu de test
24 X_test=np.transpose(X_test)
25 y_test=np.transpose(y_test.reshape((y_test.shape[0], 1)))
```

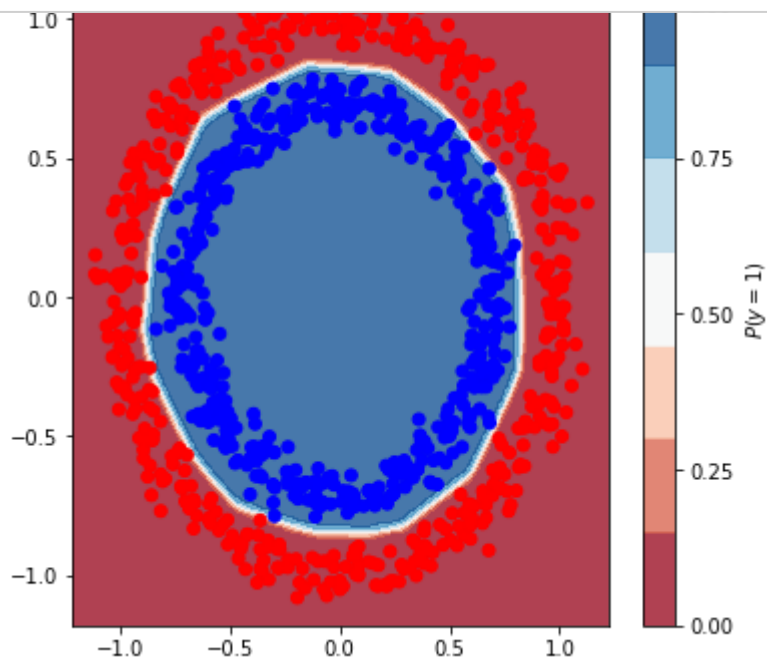
Création du réseau en utilisant *momentum*.

In [55]:

```

1  Myoptimizer="momentum"
2  Myepochs = 100
3  Myeta = 0.01
4  Mybatchsize=32
5  Mybeta=0.9
6
7  network = MyNeuralNetwork(optimizer=Myoptimizer)
8
9  network.addLayer(Layer(25,input=2,activation="relu"))
10 network.addLayer(Layer(25,activation="relu"))
11 network.addLayer(Layer(3,activation="relu"))
12 network.addLayer(Layer(1,activation="sigmoid"))
13
14
15
16 #Entraînement du classifieur
17 layers,cost_history,accuracy_history,parameter_history=network.fit(X_train, y
18                                                                    eta=Myeta, batchsize=Mybatch
19
20
21 #Prédiction
22 y_pred=network.predict(X_test)
23 accuracy_test = network.accuracy(y_pred, y_test)
24 print("Accuracy test: %.3f"%accuracy_test)
25
26 # Affichage des historiques
27
28 plot_histories (Myeta,Myepochs,cost_history,accuracy_history)
29 # Affichage de la frontière de décision
30 # Affichage de la frontière de décision
31 plot_decision_boundary(lambda x: network.predict(np.transpose(x)), X, y)
32
33

```



Le même en utilisant *adam*.

In [56]:

```

1  Myoptimizer="adam"
2  Myepochs = 100
3  Myeta = 0.01
4  Mybatchsize=32
5  Mybeta=0.9
6  Mybeta2=0.999
7  Myepsilon=1e-8
8
9  network = MyNeuralNetwork(optimizer=Myoptimizer)
10
11 network.addLayer(Layer(25,input=2,activation="relu"))
12 network.addLayer(Layer(25,activation="relu"))
13 network.addLayer(Layer(3,activation="relu"))
14 network.addLayer(Layer(1,activation="sigmoid"))
15
16
17
18 #Entraînement du classifieur
19 ▼ layers,cost_history,accuracy_history,parameter_history=network.fit(X_train, y
20 ▼                                     eta=Myeta, batchsize=Mybatch
21                                     epsilon=Myepsilon)
22
23
24 #Prédiction
25 y_pred=network.predict(X_test)
26 accuracy_test = network.accuracy(y_pred, y_test)
27 print("Accuracy test: %.3f"%accuracy_test)
28
29 # Affichage des historiques
30
31 plot_histories (Myeta,Myepochs,cost_history,accuracy_history)
32 # Affichage de la frontière de décision
33
34 plot_decision_boundary(lambda x: network.predict(np.transpose(x)), X, y)

```

Utilisation de softmax et d'adam sur le jeu de données IRIS.



In [57]:

```

1 url="https://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data
2 names = ['SepalLengthCm', 'SepalWidthCm',
3          'PetalLengthCm', 'PetalWidthCm',
4          'Species']
5
6 df = pd.read_csv(url, names=names)
7
8 # mélange des données
9 df=df.sample(frac=1).reset_index(drop=True)
10
11
12 array = df.values #nécessité de convertir le dataframe en numpy
13 #X matrice de variables prédictives - attention forcer le type à float
14 X = array[:,0:4].astype('float32')
15 #y vecteur de variable à prédire
16 y = array[:,4]
17
18 # normalisation de X
19 sc_X=StandardScaler()
20 X=sc_X.fit_transform(X)
21 # Conversion de la variable à prédire via OneHotEncoder
22 # Dans IRIS il y a 3 classes -> création de 3 colonnes pour y
23 # 1 colonne correspond à 1 classe -> 1 si la ligne est du type de la classe
24 # 0 sinon
25
26 # Integer encode
27 label_encoder = LabelEncoder()
28 integer_encoded = label_encoder.fit_transform(y)
29
30 # binary encode
31 onehot_encoder = OneHotEncoder(sparse=False, categories='auto')
32 integer_encoded = integer_encoded.reshape(len(integer_encoded), 1)
33 y = onehot_encoder.fit_transform(integer_encoded)
34 # Jeu de test/apprentissage
35 validation_size=0.6 #40% du jeu de données pour le test
36
37 testsize= 1-validation_size
38 seed=30
39 # séparation jeu d'apprentissage et jeu de test
40 X_train,X_test,y_train,y_test=train_test_split(X,
41                                                  y,
42                                                  train_size=validation_size,
43                                                  random_state=seed,
44                                                  test_size=testsize)
45
46
47 #transformation des données pour être au bon format
48 # X_train est de la forme : n colonnes (variables à prédire après OneHotEncod
49 # y_train est de la forme : m colonnes, n lignes (variables à prédire après O
50
51 # La transposée de X_train est de la forme : m colonnes (exemples), n lignes
52 X_train=np.transpose(X_train)
53
54 # y_train est forcé pour être un tableau à 1 ligne contenant m colonnes
55 y_train=np.transpose(y_train.reshape((y_train.shape[0], y_train.shape[1])))
56
57 # mêmes traitements pour le jeu de test
58 X_test=np.transpose(X_test)
59 y_test=np.transpose(y_test.reshape((y_test.shape[0], y_test.shape[1])))

```



```

60
61 Myoptimizer="adam"
62 Myepochs = 100
63 Myeta = 0.01
64 Mybatchsize=10
65 Mybeta=0.9
66 Mybeta2=0.999
67 Myepsilon=1e-8
68
69 network = MyNeuralNetwork(optimizer=Myoptimizer)
70
71 network.addLayer(Layer(10,input=4,activation="leakyrelu"))
72
73 network.addLayer(Layer(3,activation="softmax"))
74
75
76
77 #Entraînement du classifieur
78 ▼ layers,cost_history,accuracy_history,parameter_history=network.fit(X_train, y
79 ▼                                     eta=Myeta, batchsize=Mybatch
80                                     epsilon=Myepsilon)
81
82
83 #Prédiction
84 y_pred=network.predict(X_test)
85 accuracy_test = network.accuracy(y_pred, y_test)
86 print("Accuracy test: %.3f"%accuracy_test)
87
88 # Affichage des historiques
89
90 plot_histories (Myeta,Myepochs,cost_history,accuracy_history)
91 # Affichage de la frontière de décision
92

```

