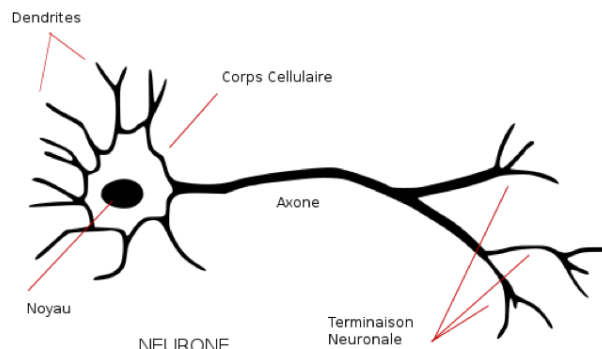


## Les réseaux de neurones

### 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)

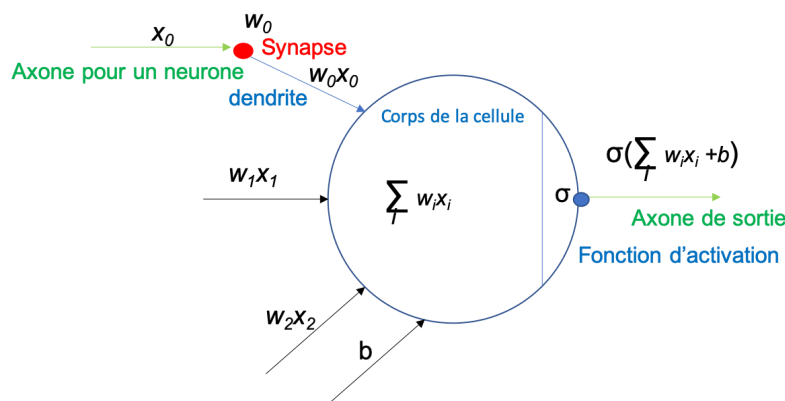
La structure d'un neurone (source : [https://fr.wikipedia.org/wiki/Fichier:Neurone\\_-\\_commenté.svg](https://fr.wikipedia.org/wiki/Fichier:Neurone_-_commenté.svg) ([https://fr.wikipedia.org/wiki/Fichier:Neurone\\_-\\_commenté.svg](https://fr.wikipedia.org/wiki/Fichier:Neurone_-_commenté.svg))) :



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 0x118b600f0>

In [2]:

```
1 #from :
2 #https://github.com/ardendertat/Applied-Deep-Learning-with-Keras/blob/master/no
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")
```

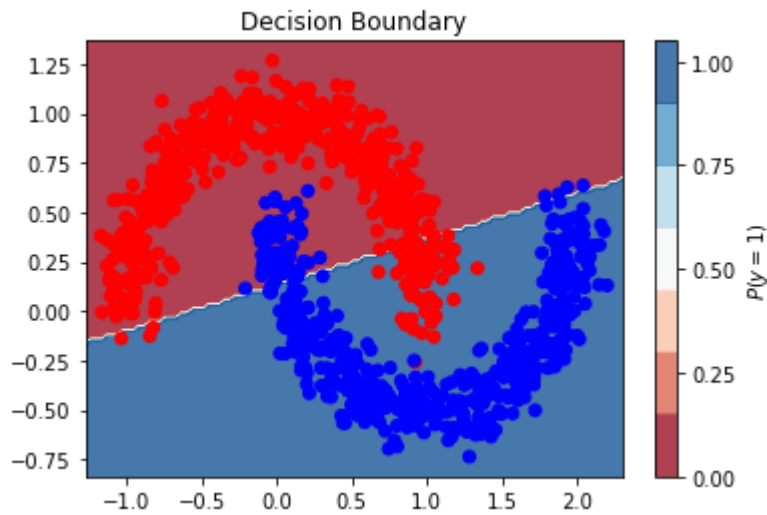
Appliquons, à présent, la logistic regression de sickit learn pour afficher la frontière de decision :

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 :

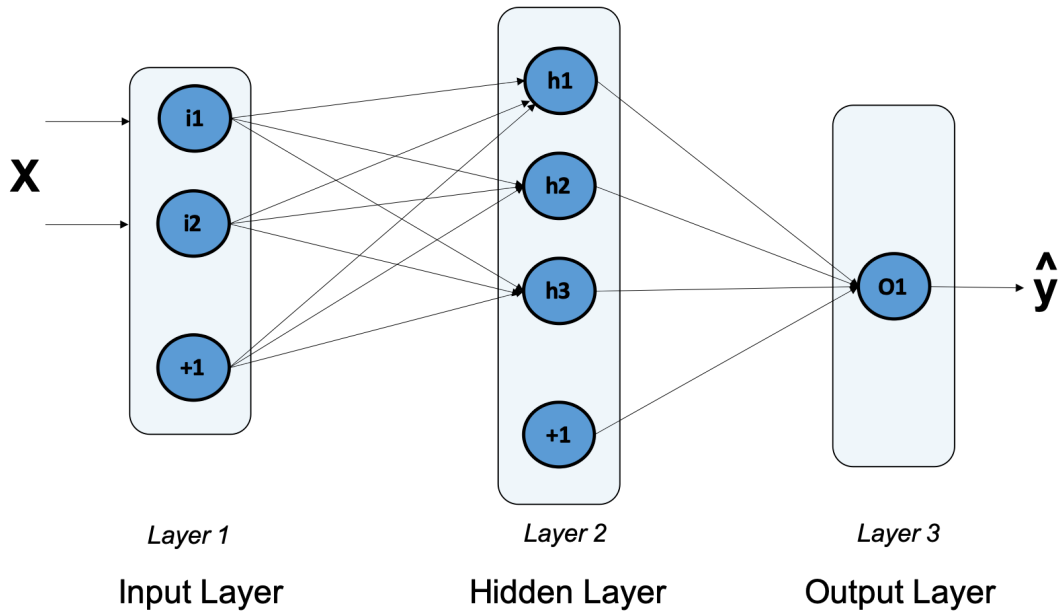
- Une couche d'entrée qui reçoit l'ensemble des caractéristiques (features), i.e. les variables prédictives.
- Un nombre arbitraire de couches cachées.
- Une couche de sortie,  $\hat{y}$ , qui contient la variable à prédire.
- Un ensemble de poids  $W$  qui vont être ajoutés aux valeurs des features 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 :

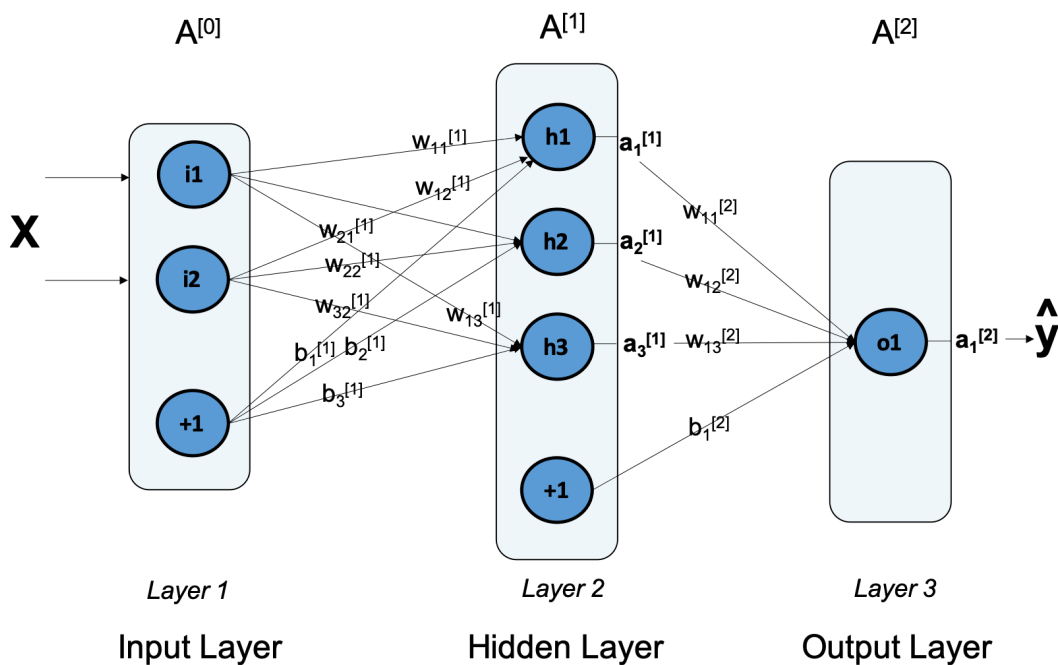
- *régression* : 1 seul neurone (C.f. notebook descente de gradient)
- *classification binaire* : 1 seul neurone avec une fonction d'activation qui sépare les deux classes.
- *classification multi-classe* : 1 neurone par classe et une fonction d'activation Softmax pour avoir la classe appropriée en fonction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

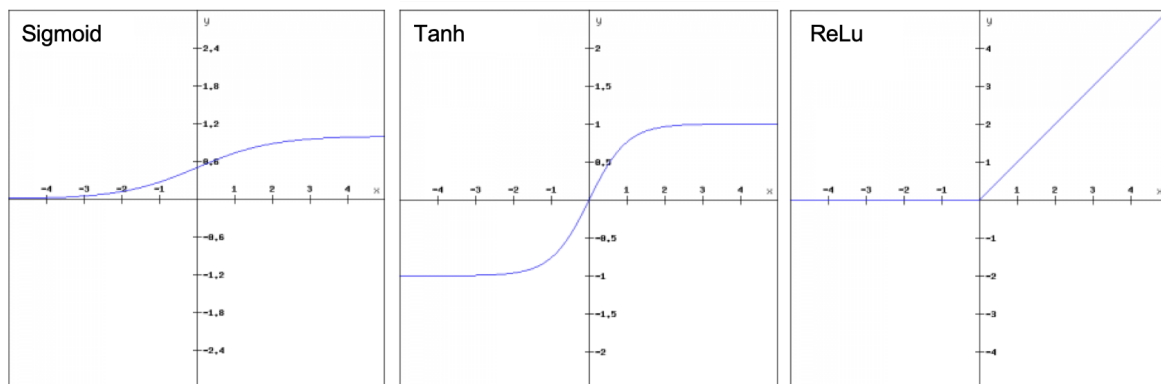


La figure suivante illustre le même réseau avec les poids affectés.



**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

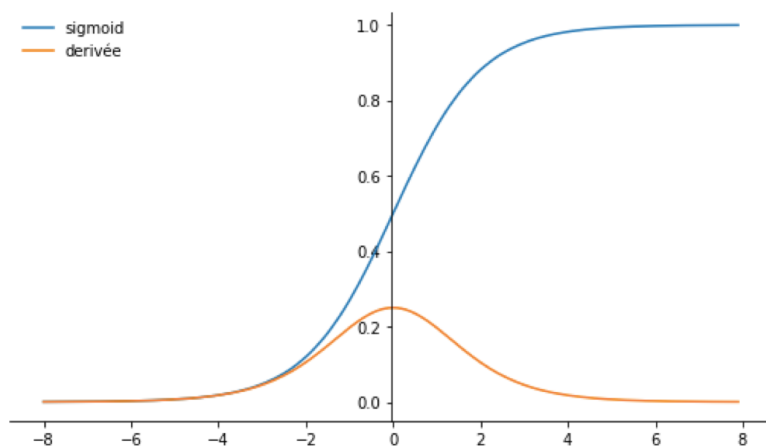


Il existe de très nombreuses fonctions d'activation qui peuvent être utilisées :

- 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 prédites :

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

Les sorties réelles 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)*

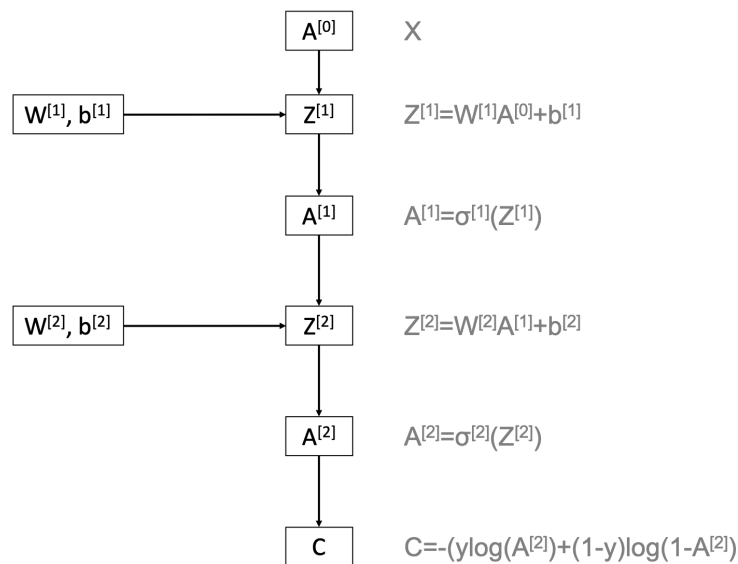
**Comprendre ce qui est derrière**

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  $\mathbf{b}^{[2]}$ .

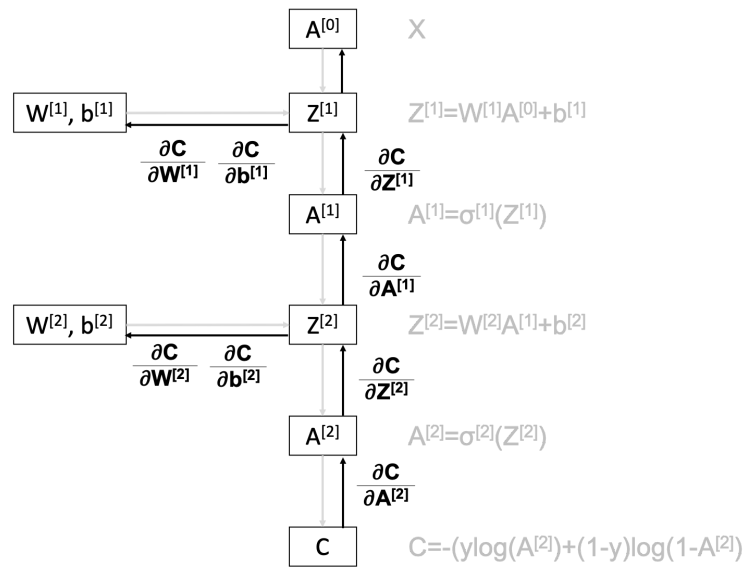
$$\begin{aligned}
 \frac{\partial \mathbf{C}}{\partial \mathbf{W}^{[2]}} &= \frac{\partial \mathbf{C}}{\partial \mathbf{A}^{[2]}} \cdot \frac{\partial \mathbf{A}^{[2]}}{\partial \mathbf{Z}^{[2]}} \cdot \frac{\partial \mathbf{Z}^{[2]}}{\partial \mathbf{W}^{[2]}} \\
 \frac{\partial \mathbf{C}}{\partial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

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  $\mathbf{l}$  sont :

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

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

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

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

Par la suite, et par simplification, nous considérons que les deux fonctions d'activation dans notre réseau sont Relu (pour  $\mathbf{A}^{[1]}$ ) et sigmoid (pour  $\mathbf{A}^{[2]}$ ). 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  $\mathbf{C}$  par rapport à  $\mathbf{A}^{[L]}$  :**

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

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

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

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

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

Pour la partie droite  $-(1 - y) \log(1 - \mathbf{A}^{[L]})$ , 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  $\mathbf{L}$ , 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  $\mathbf{l}$ , 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  $\mathbf{l}$  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

- 1 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.
- 2
- 3 Pour  $i=1 \dots C$  :

```

4  $$\mathbf{a}_i=\frac{e^{z_i}}{\sum_{k=1}^C e^{z_k}}$$
5  $$ avec\ \sum_{i=1}^C=1$$
6
7  où $C$ correspond au nombre de classes.
8

```

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 indentifie 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)))

```

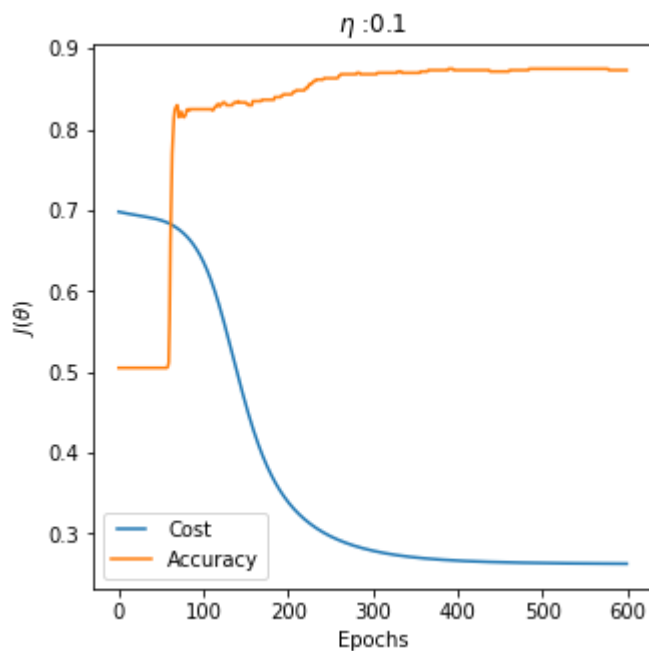
In [25]:

```

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.

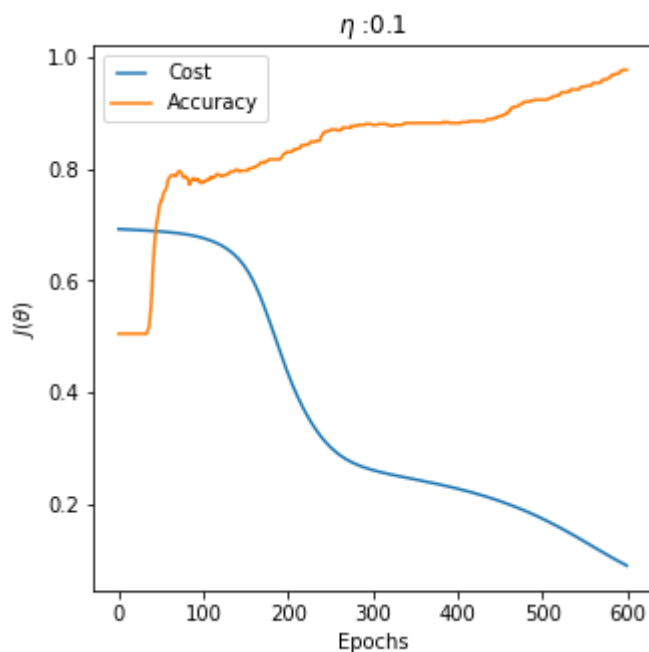
In [26]:

```

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

[...]

### Classe avec optimisation (momentum, adam) pour faire de la classification multi-classes

Dans cette section, nous présentons une classe plus complète qui permet de faire de la classification multi-classes, i.e. via la fonction d'activation softmax. En outre, elle propose de pouvoir utiliser des optimisations : momentum et adam.

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

La fonction predict est changée pour pouvoir prédire la classe d'appartenance lors de multi-classes.

La fonction update prend en compte le fait que l'optimisation peut être de type momentum ou adam.

Le constructeur prend par défaut *optimizer="bgd"* (*batch gradient descent*) (les valeurs possibles sont "*momentum*" ou "*adam*". Enfin lors de l'appel de la fonction fit les paramètres beta, beta1 et epsilon doivent être spécifiées (en cas d'appel d'un optimizer. Ici il y a des valeurs par défaut : beta=0.9, beta2=0.999, epsilon=1e-8.

La construction du réseau est similaire à celui de la classe précédente :

Pour l'exemple du notebook :

```

layers = [
{"input_dim": 2, "output_dim": 3, "activation": "relu"},

```

```
{"input_dim": 3, "output_dim": 1, "activation": "sigmoid"}  
]
```

In [39]:

```
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
```

In [40]:

```
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)
```



In [41]:

```

1  class MyNeuralNetwork(object):
2
3  def __init__(self, layers,optimizer="bgd"):
4      seed=30
5      np.random.seed(seed)
6      self.L = len(layers) #L est la couche de sortie du réseau
7      self.parameters = {}
8      self.derivatives = {}
9      self.optimizer=optimizer
10     if optimizer=="momentum":
11         self.v={} # vitesse pour momentum
12     elif optimizer == "adam":
13         self.v={}
14         self.s={}
15         self.t=2
16
17
18     # Pour toutes les couches du réseau
19     for idx, layer in enumerate(layers):
20
21         # Par simplification on commence la numérotation du layer à 1
22         # correspond aux données d'entrées (input), i.e. les variables p
23         layer_idx = idx + 1
24
25         layer_input_size = layer["input_dim"]
26         layer_output_size = layer["output_dim"]
27
28         # Initialisation des valeurs de la matrice W et du vecteur b
29         # pour les différentes couches
30         # W est initialisé en prenant l'optimisation de He et al 2015
31         self.parameters['W' + str(layer_idx)] = np.random.randn(
32             layer_output_size, layer_input_size) * np.sqrt(2/layer_input_
33         self.parameters['b' + str(layer_idx)] = np.random.randn(
34             layer_output_size, 1) * 0.1
35
36         # Sauvegarde de la fonction d'activation du réseau
37         self.parameters['activation'+ str(layer_idx)]=layer["activation"]
38         if optimizer=="momentum":
39             # Sauvegarde velocity, v, pour momentum
40             self.v['dw' + str(layer_idx)] = np.zeros_like(self.parameters
41             self.v['db' + str(layer_idx)] = np.zeros_like(self.parameters
42         elif optimizer=="adam":
43             self.v['dw' + str(layer_idx)] = np.zeros_like(self.parameters
44             self.v['db' + str(layer_idx)] = np.zeros_like(self.parameters
45             self.s['dw' + str(layer_idx)] = np.zeros_like(self.parameters
46             self.s['db' + str(layer_idx)] = np.zeros_like(self.parameters
47
48
49     def forward_propagation(self, X):
50         #Initialisation des variables prédictives pour la couche d'entrée
51         self.parameters['A0'] = X
52
53         #Propagation pour tous les layers
54         for l in range(1, self.L + 1):
55
56             # Calcul de Z
57             self.parameters['Z' + str(l)] = np.dot(self.parameters['W' + str
58                 self.parameters['A' + str(l - 1)]
59

```

```

60         # Récupération de la fonction d'activation du layer
61         activ_function_curr = self.parameters['activation' + str(l)]
62         if activ_function_curr == "relu":
63             activation_func = relu
64         elif activ_function_curr == "sigmoid":
65             activation_func = sigmoid
66         elif activ_function_curr == "tanh":
67             activation_func = tanh
68         elif activ_function_curr == "leakyrelu":
69             activation_func = leakyrelu
70         elif activ_function_curr == "softmax":
71             activation_func = softmax
72
73         # Application de la fonction d'activation à Z
74         self.parameters['A' + str(l)] = activation_func(self.parameters[
75
76
77     def convert_prob_into_class(self, probs):
78         probs = np.copy(probs) #pour ne pas perdre probs, i.e. y_hat
79         probs[probs > 0.5] = 1
80         probs[probs <= 0.5] = 0
81         return probs
82
83
84     def accuracy(self, y_hat, y):
85         if self.parameters['activation' + str(self.L)] == "softmax":
86             # si la fonction est softmax, les valeurs sont sur différentes d
87             # il faut utiliser argmax avec axis=0 pour avoir un vecteur qui
88             # où est la valeur maximale à la fois pour y_hat et pour y
89             # comme cela il suffit de comparer les deux vecteurs qui indiquen
90             # dans quelle ligne se trouve le max
91             y_hat_encoded = np.copy(y_hat)
92             y_hat_encoded = np.argmax(y_hat_encoded, axis=0)
93             y_encoded = np.copy(y)
94             y_encoded = np.argmax(y_encoded, axis=0)
95             return (y_hat_encoded == y_encoded).mean()
96         # la dernière fonction d'activation n'est pas softmax.
97         # par exemple sigmoid pour une classification binaire
98         # il suffit de convertir la probabilité du résultat en classe
99         y_hat_ = self.convert_prob_into_class(y_hat)
100        return (y_hat_ == y).all(axis=0).mean()
101
102
103     def cost_function(self, y):
104         # cross entropy
105         return -(y*np.log(self.parameters['A' + str(self.L)]+1e-8) + (1-y)*r
106
107
108     def backward_propagation(self, y):
109         #Dérivées partielles de la fonction de coût pour z[L], W[L] et b[L]:
110         #dzL
111         self.derivatives['dz' + str(self.L)] = self.parameters['A' + str(self
112         #dWL
113         self.derivatives['dW' + str(self.L)] = np.dot(self.derivatives['dz' +
114                                                         np.transpose(self.parame
115         #dbL
116         #Attention pour un lot de m exemples, il faut faire la moyenne des d
117         m = self.parameters['A' + str(self.L)].shape[1]
118         self.derivatives['db' + str(self.L)] = np.sum(self.derivatives['dz' +
119                                                         axis=1, keepdims=True)
120

```

```

121     #Dérivées partielles de la fonction de coût pour z[l], W[l] et b[l]
122     for l in range(self.L-1, 0, -1):
123
124         #Récupération de la dérivée de la fonction d'activation
125         activ_function_backward = self.parameters['activation' + str(l)]
126         if activ_function_backward == "relu":
127             backward_activation_func = relu_prime
128         elif activ_function_backward == "sigmoid":
129             backward_activation_func = sigmoid_prime
130         elif activ_function_backward == "tanh":
131             backward_activation_func = tanh_prime
132         elif activ_function_backward == "leakyrelu":
133             backward_activation_func = leakyrelu_prime
134
135         self.derivatives['dz' + str(l)] = np.dot(np.transpose(self.parameters['W' + str(l + 1)]),
136             self.derivatives['dz' + str(l + 1)]) * backward_activation_func
137
138         self.derivatives['dw' + str(l)] = np.dot(self.derivatives['dz' + str(l)],
139             np.transpose(self.parameters['A' + str(l - 1)]))
140         #Attention pour un lot de m exemples, il faut faire la moyenne de
141         m = self.parameters['A' + str(l - 1)].shape[1]
142         self.derivatives['db' + str(l)] = np.sum(self.derivatives['dz' + str(l)],
143             axis=1, keepdims=True)
144
145     def update_parameters(self, eta, beta, beta2=0.999, epsilon=1e-8):
146         # Descente de gradient
147         if self.optimizer=="adam":
148             v_corrected = {} # Initialisation de la première estimation
149             s_corrected = {} # Initialisation de la seconde estimation
150
151         for l in range(1, self.L+1):
152             if self.optimizer=="momentum":
153                 # Calcul de la vitesse
154                 self.v['dw' + str(l)] = beta * self.v['dw' + str(l)] + (1 - beta) * self.derivatives['dw' + str(l)]
155                 self.v['db' + str(l)] = beta * self.v['db' + str(l)] + (1 - beta) * self.derivatives['db' + str(l)]
156
157                 # Mise à jour des paramètres
158                 self.parameters['W' + str(l)] -= eta * self.v['dw' + str(l)]
159                 self.parameters['b' + str(l)] -= eta * self.v['db' + str(l)]
160             elif self.optimizer=="adam":
161                 # Calcul de la vitesse
162                 self.v['dw' + str(l)] = beta * self.v['dw' + str(l)] + (1 - beta) * self.derivatives['dw' + str(l)]
163                 self.v['db' + str(l)] = beta * self.v['db' + str(l)] + (1 - beta) * self.derivatives['db' + str(l)]
164
165                 # Calcul de la première estimation du moment (correction du biais)
166                 v_corrected["dw" + str(l)] = self.v["dw" + str(l)] / (1 - np.power(beta, l))
167                 v_corrected["db" + str(l)] = self.v["db" + str(l)] / (1 - np.power(beta, l))
168
169                 # déplacement moyen des gradients au carré
170                 self.s["dw" + str(l)] = beta2 * self.s["dw" + str(l)] + (1 - beta2) * v_corrected["dw" + str(l)] ** 2
171                 self.s["db" + str(l)] = beta2 * self.s["db" + str(l)] + (1 - beta2) * v_corrected["db" + str(l)] ** 2
172
173                 # Calcul de la seconde estimation du moment (correction du biais)
174                 s_corrected["dw" + str(l)] = self.s["dw" + str(l)] / (1 - np.power(beta2, l))
175                 s_corrected["db" + str(l)] = self.s["db" + str(l)] / (1 - np.power(beta2, l))
176
177                 # Mise à jour des paramètres
178                 self.parameters['W' + str(l)] -= eta * v_corrected["dw" + str(l)]
179                 self.parameters['b' + str(l)] -= eta * v_corrected["db" + str(l)]
180
181             else: #descente par mini-lots

```

```

182         self.parameters['W' + str(l)] -= eta*self.derivatives['dW' +
183         self.parameters['b' + str(l)] -= eta*self.derivatives['db' +
184
185
186     def predict(self, x):
187         self.forward_propagation(x)
188         return self.parameters['A' + str(self.L)]
189
190
191     def next_batch(self,X, y, batchsize):
192         # pour avoir X de la forme : 2 colonnes, m lignes (exemples) et égale
193         # cela permet de trier les 2 tableaux avec un indices de permutation
194         X=np.transpose(X)
195         y=np.transpose(y)
196
197         m=len(y)
198         # permutation aléatoire de X et y pour faire des batchs avec des valeurs
199         indices = np.random.permutation(m)
200         X = X[indices]
201         y = y[indices]
202         for i in np.arange(0, X.shape[0], batchsize):
203             # creation des batchs de taille batchsize
204             yield (X[i:i + batchsize], y[i:i + batchsize])
205
206
207
208     def fit(self, X, y, epochs, eta = 0.01,batchsize=32,beta=0.9,beta2=0.999,
209
210         # sauvegarde historique coût et accuracy pour affichage
211         cost_history = []
212         accuracy_history = []
213         for i in range(epochs):
214             # sauvegarde des coûts et accuracy par mini-batch
215             cost_batch = []
216             accuracy_batch = []
217             # Descente de gradient par mini-batch
218             for (batchX, batchy) in self.next_batch(X, y, batchsize):
219                 # Extraction et traitement d'un batch à la fois
220
221                 # mise en place des données au bon format
222                 batchX=np.transpose(batchX)
223                 if self.parameters['activation' + str(self.L)]=="softmax":
224                     # la classification n'est pas binaire, y a utilisé one-hot
225                     # le batchy doit donc être transposé et le résultat doit
226                     # être sous la forme d'une matrice de taille batchy.shape[0], batchy.shape[1]
227                     batchy=np.transpose(batchy.reshape((batchy.shape[0], batchy.shape[1])))
228                 else:
229                     # il s'agit d'une classification binaire donc shape[1] n'est pas 2
230                     # pas
231                     batchy=np.transpose(batchy.reshape((batchy.shape[0], 1)))
232
233
234                 self.forward_propagation(batchX)
235                 self.backward_propagation(batchy)
236                 if self.optimizer=="adam":
237                     self.t=self.t+1
238                     self.update_parameters(eta,beta,beta2,epsilon)
239                 else:
240                     # self.update_parameters(eta,0)
241
242                 # sauvegarde pour affichage

```

```

243         current_cost=self.cost_function(batchy)
244         cost_batch.append(current_cost)
245         y_hat = self.predict(batchX)
246         current_accuracy = self.accuracy(y_hat, batchy)
247         accuracy_batch.append(current_accuracy)
248
249
250         # sauvegarde de la valeur moyenne des coûts et de l'accuracy du lot
251         current_cost=np.average(cost_batch)
252         cost_history.append(current_cost)
253         current_accuracy=np.average(accuracy_batch)
254         accuracy_history.append(current_accuracy)
255
256     if(i % 10 == 0):
257         print("Epoch : #%-s - cost : %.3f - accuracy : %.3f"%(i, float
258         return self.parameters, cost_history, accuracy_history

```

## Essai classe avec optimisation sur différents jeux de données - classification binaire

Sélection du jeu de données

In [42]:

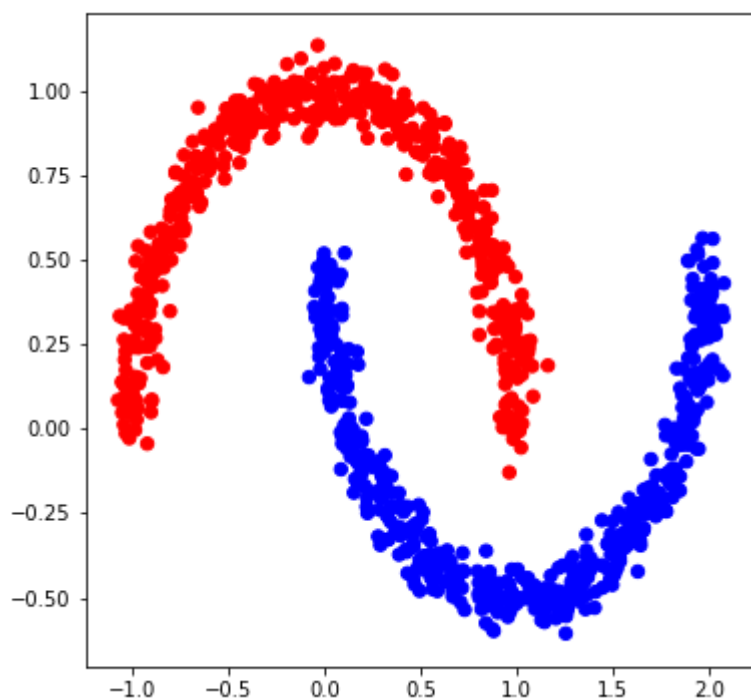
```

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[42]:

<matplotlib.collections.PathCollection at 0x12f324a90>



Le même en utilisant l'optimisation Momentum.

In [43]:

```

1  print ("Avec Momentum")
2  #Création du classifieur (mise en place des layers et initialisation des W et
3  optimizer="momentum"
4  clf = MyNeuralNetwork(layers, optimizer)
5  beta=0.9
6  #Entraînement du classifieur
7  parameters,cost_history,accuracy_history=clf.fit(X_train, y_train, epochs,eta
8
9  y_pred=clf.predict(X_test)
10 accuracy_test = clf.accuracy(y_pred, y_test)
11 print("Accuracy : %.3f"%accuracy_test)
12 # Affichage des historiques
13 plot_histories (eta,epochs,cost_history,accuracy_history)
14 # Affichage de la frontière de décision
15 plot_decision_boundary(lambda x: clf.predict(np.transpose(x)), X, y)

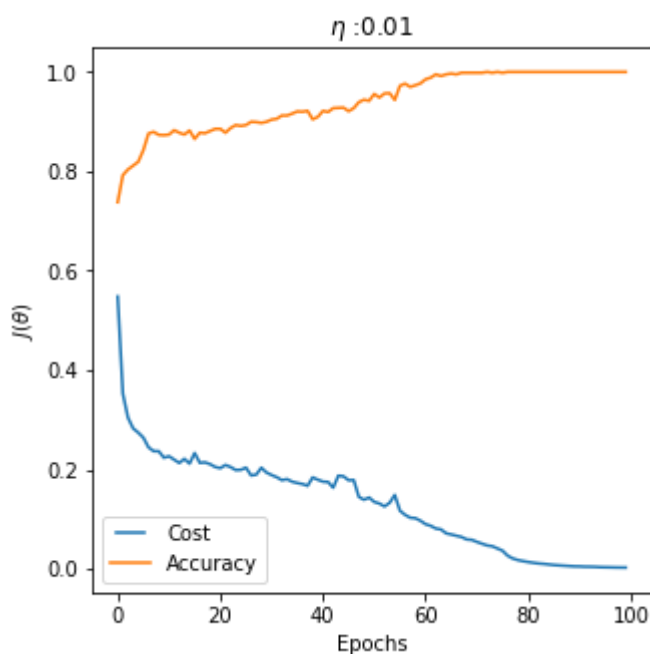
```

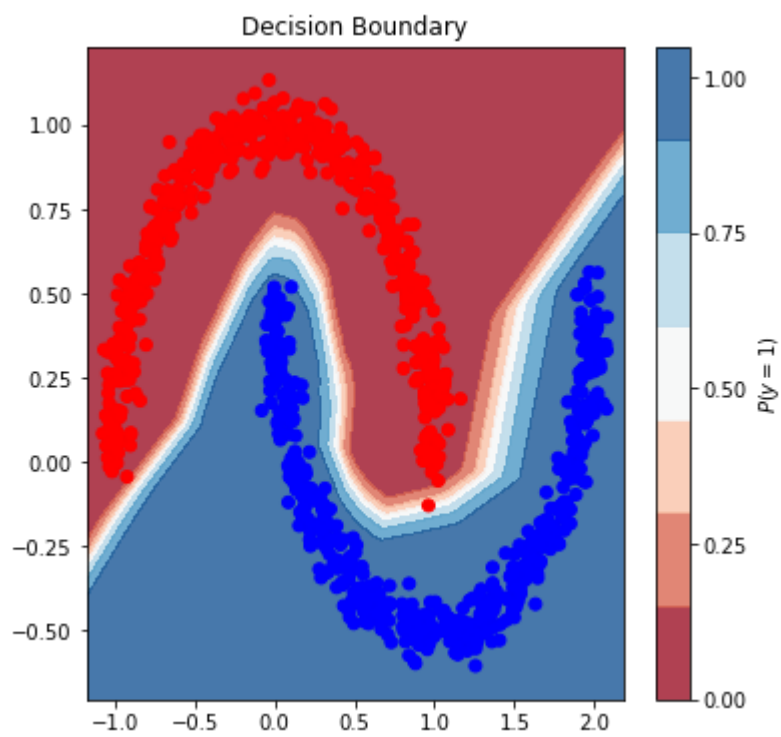
Avec Momentum

```

Epoch : #0 - cost : 0.548 - accuracy : 0.738
Epoch : #10 - cost : 0.227 - accuracy : 0.874
Epoch : #20 - cost : 0.203 - accuracy : 0.885
Epoch : #30 - cost : 0.189 - accuracy : 0.904
Epoch : #40 - cost : 0.176 - accuracy : 0.922
Epoch : #50 - cost : 0.135 - accuracy : 0.956
Epoch : #60 - cost : 0.090 - accuracy : 0.985
Epoch : #70 - cost : 0.054 - accuracy : 0.998
Epoch : #80 - cost : 0.014 - accuracy : 1.000
Epoch : #90 - cost : 0.005 - accuracy : 1.000
Accuracy : 1.000

```





Le même avec adam.

In [44]:

```

1  print ("Avec Adam")
2  #Création du classifieur (mise en place des layers et initialisation des W et
3  optimizer="adam"
4  clf = MyNeuralNetwork(layers, optimizer)
5  beta=0.9
6  beta2=0.999
7  epsilon=1e-8
8  #Entraînement du classifieur
9  parameters,cost_history,accuracy_history=clf.fit(X_train, y_train, epochs,eta
10
11  y_pred=clf.predict(X_test)
12  accuracy_test = clf.accuracy(y_pred, y_test)
13  print("Accuracy test : %.3f"%accuracy_test)
14  # Affichage des historiques
15  plot_histories (eta,epochs,cost_history,accuracy_history)
16  # Affichage de la frontière de décision
17  plot_decision_boundary(lambda x: clf.predict(np.transpose(x)), X, y)

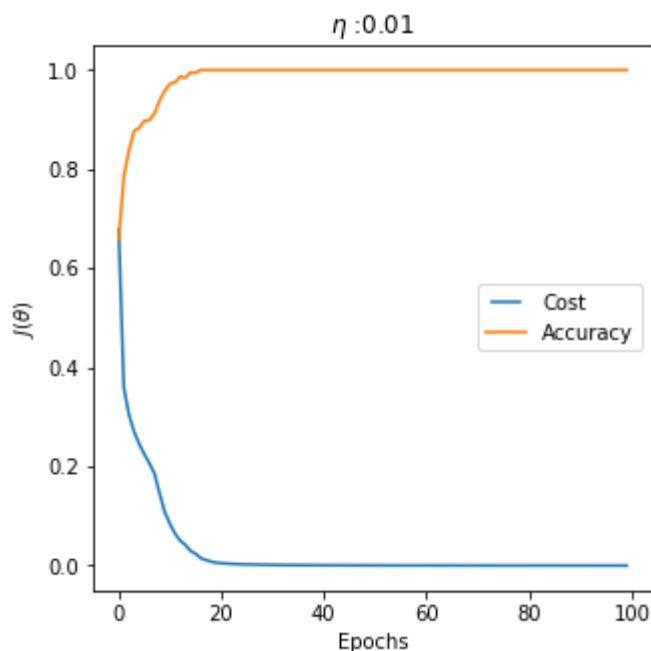
```

Avec Adam

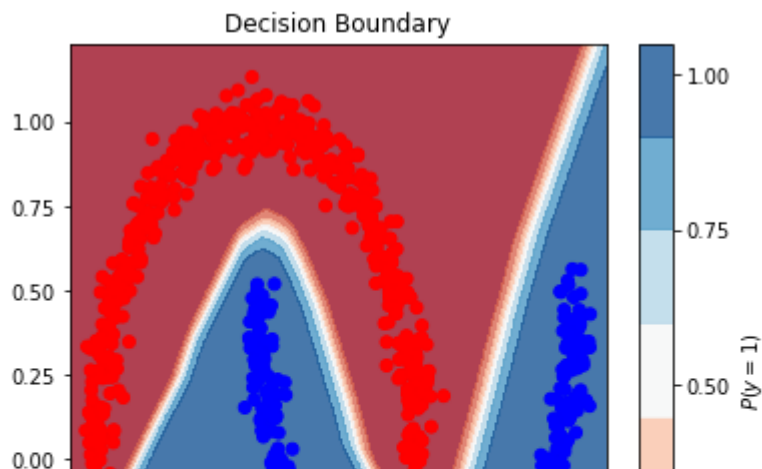
```

Epoch : #0 - cost : 0.678 - accuracy : 0.659
Epoch : #10 - cost : 0.084 - accuracy : 0.972
Epoch : #20 - cost : 0.005 - accuracy : 1.000
Epoch : #30 - cost : 0.001 - accuracy : 1.000
Epoch : #40 - cost : 0.001 - accuracy : 1.000
Epoch : #50 - cost : 0.000 - accuracy : 1.000
Epoch : #60 - cost : 0.000 - accuracy : 1.000
Epoch : #70 - cost : 0.000 - accuracy : 1.000
Epoch : #80 - cost : 0.000 - accuracy : 1.000
Epoch : #90 - cost : 0.000 - accuracy : 1.000
Accuracy test : 1.000

```







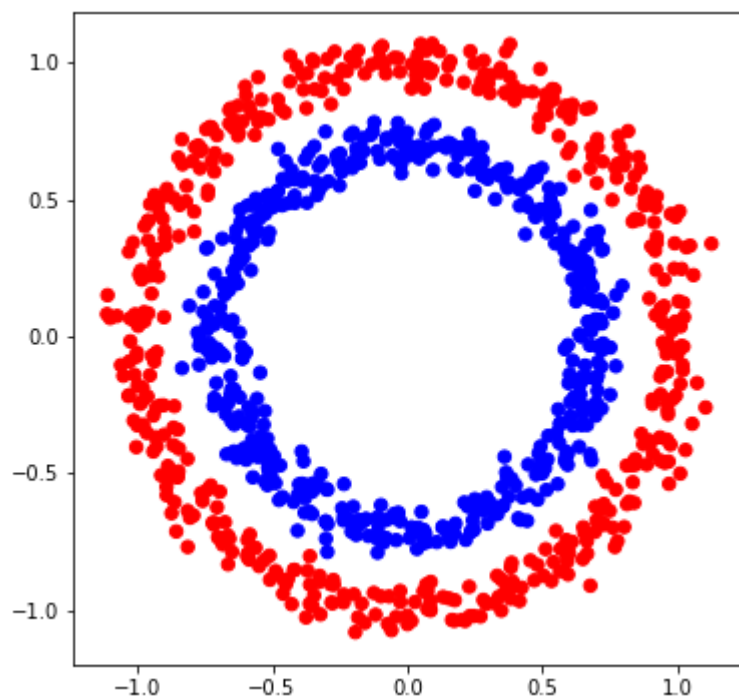
Test sur un autre jeu de données.

In [45]:

```
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[45]:

<matplotlib.collections.PathCollection at 0x12e43db00>



In [46]:

```

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)))
26
27
28 # déclaration de l'architecture du réseau
29 layers = [
30     {"input_dim": 2, "output_dim": 25, "activation": "relu"},
31     {"input_dim": 25, "output_dim": 3, "activation": "relu"},
32     {"input_dim": 3, "output_dim": 1, "activation": "sigmoid"},
33 ]
34
35
36 epochs = 200
37 eta = 0.01
38 beta=0.9
39 beta2=0.999
40 epsilon=1e-8
41 #Création du classifieur (mise en place des layers et initialisation des W et
42 print ("Batch gradient descent")
43 clf = MyNeuralNetwork(layers)
44 #Entraînement du classifieur
45 parameters,cost_history,accuracy_history=clf.fit(X_train, y_train, epochs,eta
46 plot_histories (eta,epochs,cost_history,accuracy_history)
47 print ("\nBatch gradient descent et momentum")
48 clf = MyNeuralNetwork(layers,optimizer="momentum")
49 #Entraînement du classifieur
50 parameters,cost_history,accuracy_history=clf.fit(X_train, y_train, epochs,eta
51 plot_histories (eta,epochs,cost_history,accuracy_history)
52 print ("\nBatch gradient descent et adam")
53 clf = MyNeuralNetwork(layers,optimizer="momentum")
54 #Entraînement du classifieur
55 parameters,cost_history,accuracy_history=clf.fit(X_train, y_train, epochs,eta
56 plot_histories (eta,epochs,cost_history,accuracy_history)
57
58 #Prédiction
59 y_pred= clf.predict(X_test)

```

```

60 accuracy_test = clf.accuracy(y_pred, y_test)
61 print("Accuracy test : %.3f"%accuracy_test)
62
63 # Affichage de la frontière de décision
64 plot_decision_boundary(lambda x: clf.predict(np.transpose(x)), X, y)

```

#### Batch gradient descent

```

Epoch : #0 - cost : 0.709 - accuracy : 0.522
Epoch : #10 - cost : 0.638 - accuracy : 0.621
Epoch : #20 - cost : 0.521 - accuracy : 0.817
Epoch : #30 - cost : 0.236 - accuracy : 0.949
Epoch : #40 - cost : 0.069 - accuracy : 0.996
Epoch : #50 - cost : 0.027 - accuracy : 1.000
Epoch : #60 - cost : 0.356 - accuracy : 0.940
Epoch : #70 - cost : 0.017 - accuracy : 1.000
Epoch : #80 - cost : 0.011 - accuracy : 1.000
Epoch : #90 - cost : 0.494 - accuracy : 0.947
Epoch : #100 - cost : 0.009 - accuracy : 1.000
Epoch : #110 - cost : 0.007 - accuracy : 1.000
Epoch : #120 - cost : 0.006 - accuracy : 1.000
Epoch : #130 - cost : 0.005 - accuracy : 1.000
Epoch : #140 - cost : 0.006 - accuracy : 1.000
Epoch : #150 - cost : 0.004 - accuracy : 1.000
Epoch : #160 - cost : 0.004 - accuracy : 1.000
Epoch : #170 - cost : 0.004 - accuracy : 1.000
Epoch : #180 - cost : 0.004 - accuracy : 1.000
Epoch : #190 - cost : 0.003 - accuracy : 1.000

```

#### Batch gradient descent et momentum

```

Epoch : #0 - cost : 0.742 - accuracy : 0.497
Epoch : #10 - cost : 0.636 - accuracy : 0.515
Epoch : #20 - cost : 0.449 - accuracy : 0.874
Epoch : #30 - cost : 0.113 - accuracy : 0.988
Epoch : #40 - cost : 0.038 - accuracy : 0.998
Epoch : #50 - cost : 0.015 - accuracy : 1.000
Epoch : #60 - cost : 0.012 - accuracy : 0.998
Epoch : #70 - cost : 0.009 - accuracy : 1.000
Epoch : #80 - cost : 0.025 - accuracy : 0.998
Epoch : #90 - cost : 0.005 - accuracy : 1.000
Epoch : #100 - cost : 0.040 - accuracy : 0.995
Epoch : #110 - cost : 0.038 - accuracy : 0.992
Epoch : #120 - cost : 0.032 - accuracy : 0.995
Epoch : #130 - cost : 0.005 - accuracy : 1.000
Epoch : #140 - cost : 0.002 - accuracy : 1.000
Epoch : #150 - cost : 0.003 - accuracy : 1.000
Epoch : #160 - cost : 0.004 - accuracy : 1.000
Epoch : #170 - cost : 0.002 - accuracy : 1.000
Epoch : #180 - cost : 0.002 - accuracy : 1.000
Epoch : #190 - cost : 0.009 - accuracy : 1.000

```

#### Batch gradient descent et adam

```

Epoch : #0 - cost : 0.742 - accuracy : 0.497
Epoch : #10 - cost : 0.636 - accuracy : 0.515
Epoch : #20 - cost : 0.449 - accuracy : 0.874
Epoch : #30 - cost : 0.113 - accuracy : 0.988
Epoch : #40 - cost : 0.038 - accuracy : 0.998
Epoch : #50 - cost : 0.015 - accuracy : 1.000
Epoch : #60 - cost : 0.012 - accuracy : 0.998
Epoch : #70 - cost : 0.009 - accuracy : 1.000
Epoch : #80 - cost : 0.025 - accuracy : 0.998
Epoch : #90 - cost : 0.005 - accuracy : 1.000

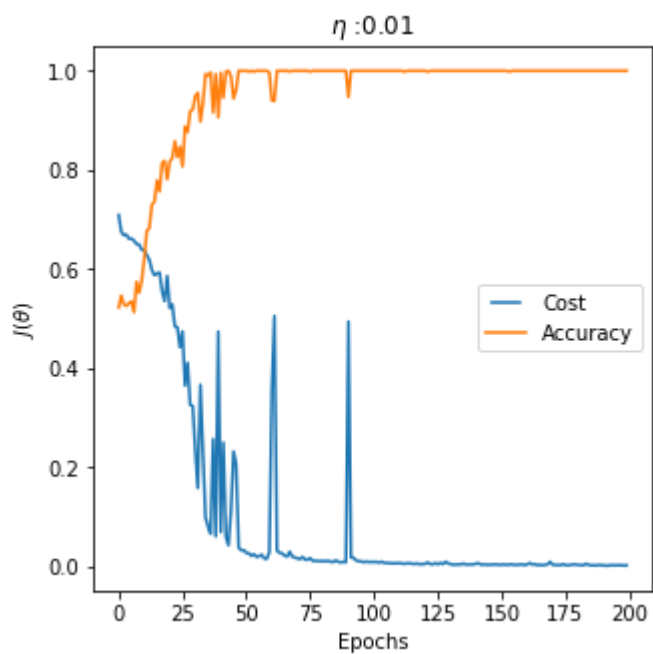
```

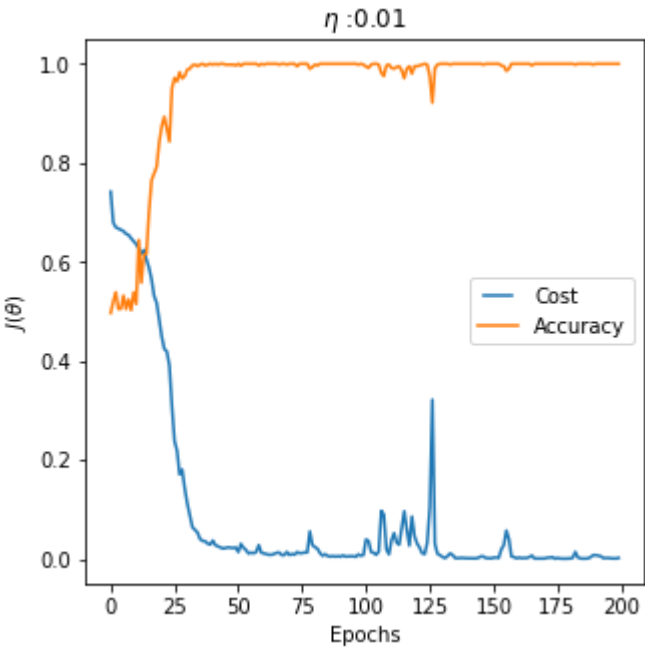
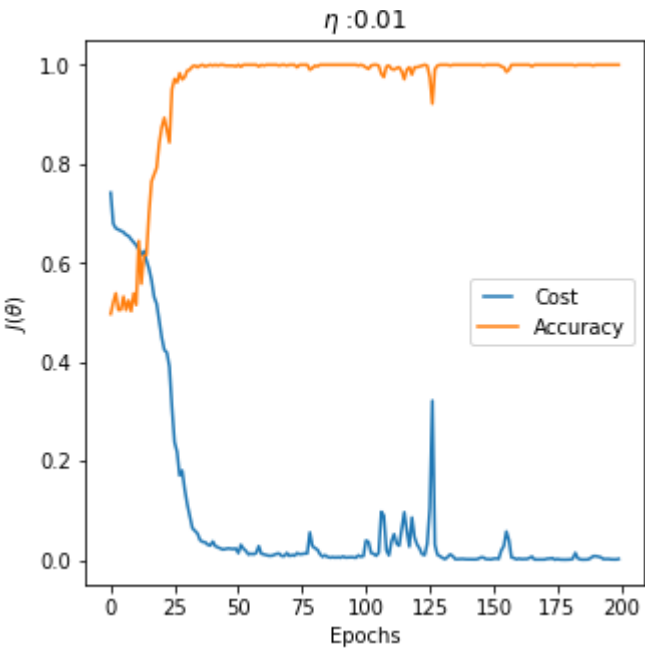
```

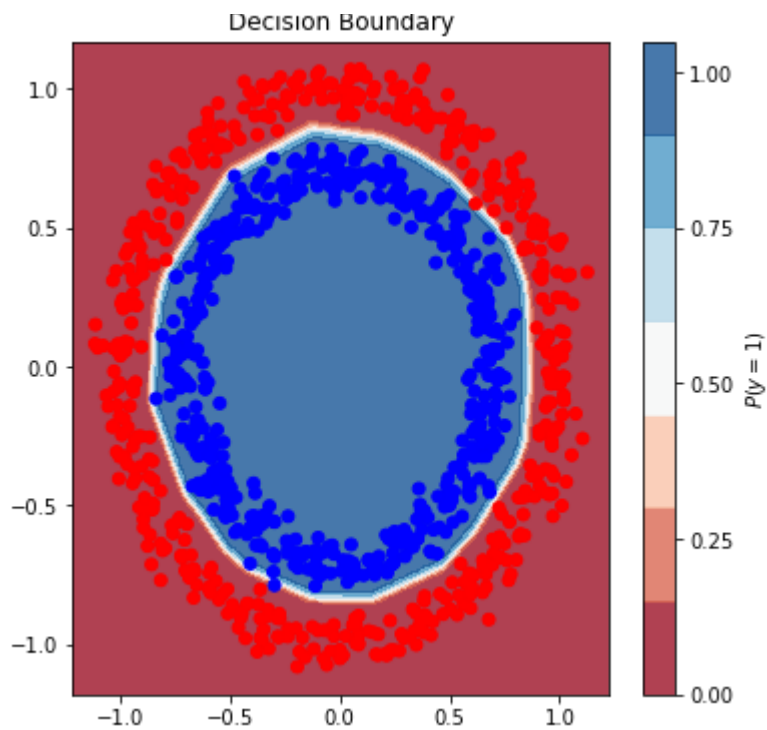
Epoch : #100 - cost : 0.040 - accuracy : 0.995

```

```
Epoch : #100 - cost : 0.040 - accuracy : 0.995  
Epoch : #110 - cost : 0.038 - accuracy : 0.992  
Epoch : #120 - cost : 0.032 - accuracy : 0.995  
Epoch : #130 - cost : 0.005 - accuracy : 1.000  
Epoch : #140 - cost : 0.002 - accuracy : 1.000  
Epoch : #150 - cost : 0.003 - accuracy : 1.000  
Epoch : #160 - cost : 0.004 - accuracy : 1.000  
Epoch : #170 - cost : 0.002 - accuracy : 1.000  
Epoch : #180 - cost : 0.002 - accuracy : 1.000  
Epoch : #190 - cost : 0.009 - accuracy : 1.000  
Accuracy test : 0.998
```







### Essai classe avec classification multi-classes (softmax)

Dans cette section, nous illustrons l'utilisation de la fonction d'activation softmax pour faire de la classification multi-classes. Nous reprenons l'exemple d'IRIS.

In [47]:

```

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
```

Première étape : normalisation des données

In [48]:

```
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 [49]:

```
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)
14
15
```

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

In [50]:

```
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])))
27
```

Création du réseau et lancement du classifieur.



In [51]:

```

1  ▾ # déclaration de l'architecture du réseau
2  ▾ layers = [
3      {"input_dim": 4, "output_dim": 10, "activation": "leakyrelu"},
4      {"input_dim": 10, "output_dim": 3, "activation": "softmax"},
5  ]
6
7
8  #Création du classifieur (mise en place des layers et initialisation des W et
9  # optimizer = bgd (descente par mini-batch) par défaut sinon optimizer="momentum"
10 clf = MyNeuralNetwork(layers,optimizer="adam")
11
12
13 epochs = 100
14 eta = 0.01
15 batchsize=10
16 beta=0.9
17 #Entraînement du classifieur
18 parameters,cost_history,accuracy_history=clf.fit(X_train, y_train, epochs,eta)
19
20 #Prédiction
21 y_pred= clf.predict(X_test)
22 accuracy_test = clf.accuracy(y_pred, y_test)
23 print("Accuracy test: %.3f"%accuracy_test)
24
25 #Affichage des historiques
26 plot_histories (eta,epochs,cost_history,accuracy_history)
27

```

```

Epoch : #0 - cost : 0.557 - accuracy : 0.500
Epoch : #10 - cost : 0.127 - accuracy : 0.944
Epoch : #20 - cost : 0.067 - accuracy : 0.978
Epoch : #30 - cost : 0.043 - accuracy : 0.989
Epoch : #40 - cost : 0.031 - accuracy : 1.000
Epoch : #50 - cost : 0.025 - accuracy : 1.000
Epoch : #60 - cost : 0.021 - accuracy : 1.000
Epoch : #70 - cost : 0.019 - accuracy : 1.000
Epoch : #80 - cost : 0.017 - accuracy : 1.000
Epoch : #90 - cost : 0.014 - accuracy : 1.000
Accuracy test: 0.967

```

$\eta : 0.01$