

## Descente de gradient

La descente de gradient (*Gradient Descent*) est un algorithme d'optimisation qui permet de rechercher le minimum d'une fonction (minimiser une fonction) qui est différentiable (dont on peut calculer la dérivée) définie sur un espace Euclidien (par exemple  $\mathbb{R}^n$ ). Ce n'est pas un concept nouveau, il a été inventé par un mathématicien Français, [Louis Augustin Cauchy](https://www.math.uni-bielefeld.de/documenta/vol-ismp/40_lemarechal-claude.pdf) ([https://www.math.uni-bielefeld.de/documenta/vol-ismp/40\\_lemarechal-claude.pdf](https://www.math.uni-bielefeld.de/documenta/vol-ismp/40_lemarechal-claude.pdf)) en 1847. Il est par contre d'un grand intérêt aujourd'hui notamment pour l'entraînement des réseaux profonds (deep learning).

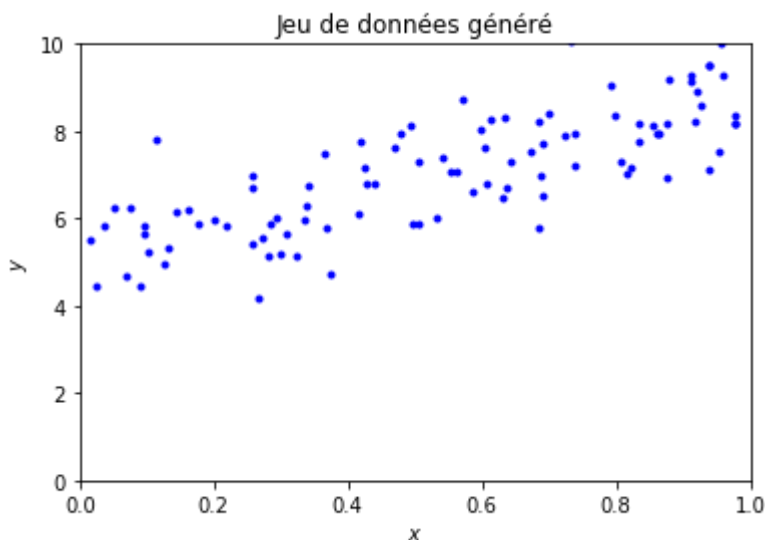
Dans la suite de ce notebook, nous nous intéressons plus particulièrement à l'apprentissage supervisé.

## Cas de la régression linéaire

De manière à mieux comprendre le concept de descente de gradient, considérons un jeu de données aléatoire pour lequel nous savons que la fonction linéaire pour le créer initialement est  $y = 4x + 5$ .

In [1]:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 %matplotlib inline
4
5 #definition de la fonction y=f(4X + 5)
6 X = np.random.rand(100,1)
7 y = 5 + 4*X+np.random.randn(100,1)
8 #X = 2 * np.random.rand(100,1)
9 #y = 4 +3 * X+np.random.randn(100,1)
10 #Visualisation du jeu de données
11 plt.plot(X,y, 'b. ')
12 plt.title('Jeu de données généré')
13 plt.xlabel("$x$", fontsize=10)
14 plt.ylabel("$y$", fontsize=10)
15 #plt.plot(X,4*X+ 5, 'r- ')
16 plt.axis([0,1,0,10])
17 plt.show()
```



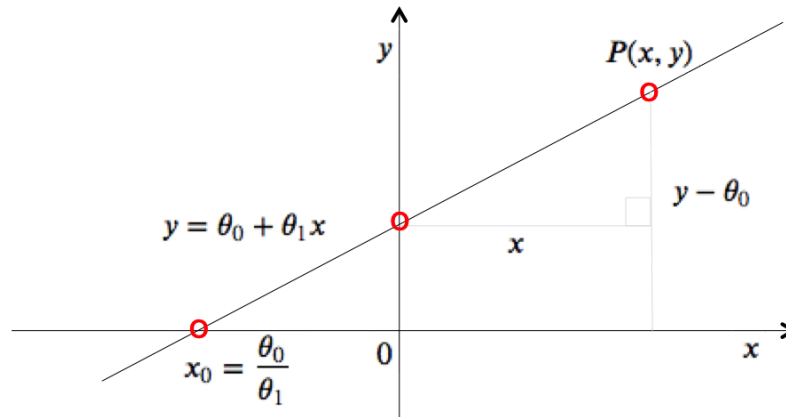
Actuellement nous connaissons plus ou moins l'équation de la droite (des valeurs aléatoires ont été ajoutées) mais l'objectif d'un algorithme d'apprentissage est, à partir d'un ensemble connu d'apprentissage de pouvoir

prédire les valeurs de nouvelles entrées. En d'autres termes, nous souhaitons construire un modèle : une hypothèse qui peut être utilisée pour estimer  $y$  en fonction de  $x$ .

Nous savons que notre hypothèse à trouver est de la forme :

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

où  $\theta_0$  correspond au point d'intersection avec l'axe des  $y$  et  $\theta_1$  à la pente de la droite.



Considérons la fonction predict suivante et initialisons  $\theta_0 = 2$  et  $\theta_1 = 1$  :

In [2]:

```
1 def predict(X, theta0, theta1):
2     return theta0 + theta1*X
3
4 print ('Valeur de X[5] : ',X[5],' - valeur réelle de y[5] : ',y[5])
5 theta0=2
6 theta1=1
7 print ('valeur prédite : ',predict(X[5], theta0, theta1))
```

```
Valeur de X[5] : [0.54168813] - valeur réelle de y[5] : [7.3796979
2]
valeur prédite : [2.54168813]
```

Nous pouvons constater que la valeur prédite est loin de la valeur réelle.

L'objectif est donc de déterminer les valeurs de  $\theta_0$  et  $\theta_1$ , de sorte que l'erreur soit la plus faible possible.

La fonction de coût (*Cost Function*) ou de perte (*Loss Function*) permet d'évaluer les erreurs de prédiction. La fonction de coût calcule l'erreur pour une seule occurrence d'apprentissage, tandis que la fonction de perte est la moyenne des fonctions de coûts pour tous les exemples d'apprentissage. Par abus de langage, les deux sont utilisés.

Si nous avons  $m$  données dans l'ensemble d'apprentissage, nous souhaitons minimiser l'erreur. Une fonction de coût souvent utilisée est l'erreur quadratique moyenne (*Mean square error (MSE)*)

[https://en.wikipedia.org/wiki/Mean\\_squared\\_error](https://en.wikipedia.org/wiki/Mean_squared_error) ([https://en.wikipedia.org/wiki/Mean\\_squared\\_error](https://en.wikipedia.org/wiki/Mean_squared_error)) :

$$J = \frac{1}{m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2$$

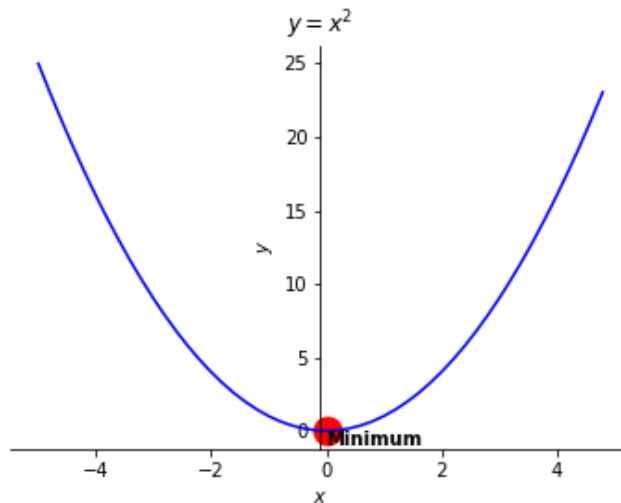
*Remarque* : les différences quadratiques sont préférées aux différences absolues car elles permettent de trouver une ligne de régression en calculant la dérivée première. En outre le fait de mettre au carré augmente la distance d'erreur ce qui met en évidence les mauvaises prédictions.

Dans notre cas,  $J$  dépend de  $\theta_0$  et  $\theta_1$  donc :

$$\begin{aligned} J(\theta_0, \theta_1) &= \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \\ &= \frac{1}{m} \sum_{i=1}^m ((\theta_0 + \theta_1 x^{(i)}) - y^{(i)})^2 \end{aligned}$$

### Comment minimiser la fonction de coût ?

En fait si nous regardons  $J$  elle est de type :  $y = x^2$ . Nous pouvons la dessiner et le minimum est situé à la position (0,0) (point rouge sur le dessin).



Pour minimiser la fonction il faut trouver la valeur de  $x$  qui produit la valeur de  $y$  la plus basse. Pour cela on utilise l'algorithme de descente de gradient.

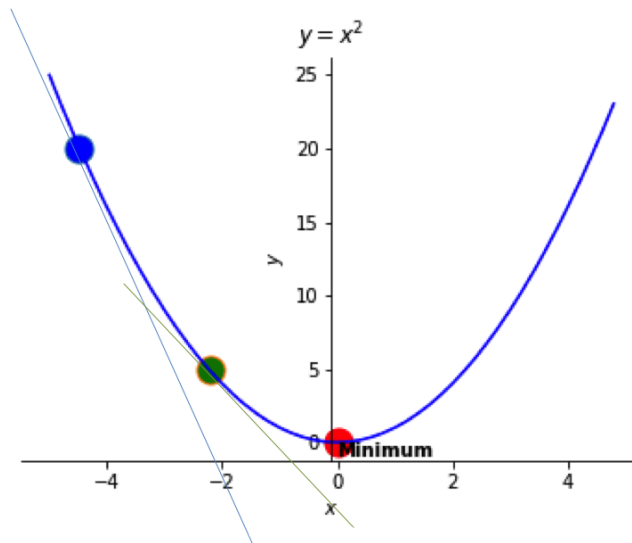
### Algorithme de Descente de Gradient (*Gradient Descent*)

La descente de gradient est un algorithme, très utilisé en apprentissage, itératif d'optimisation pour trouver le minimum d'une fonction, i.e. dans notre cas la fonction de perte.

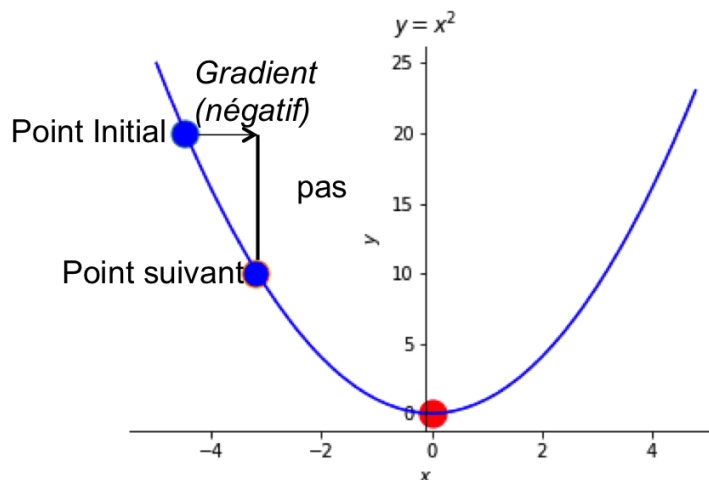
Un exemple classique pour expliquer l'approche est la suivante. Considérons une personne sur une montagne avec les yeux bandés qui souhaite descendre dans la plaine. La seule aide dont dispose la personne est l'altitude qu'elle peut relever. La personne peut commencer à descendre dans une direction aléatoire puis regarder l'altitude. Si elle est supérieure à l'altitude initiale, elle sait qu'elle est dans la mauvaise direction. Elle peut alors changer de direction et répéter le même processus. De la même manière elle peut marcher plus ou moins vite en fonction de la pente. Si la pente est raide elle descend vite et autrement plus lentement. Au bout de beaucoup d'itérations elle arrive au creux de la vallée.

Dans notre cas, les gradients indiquent la direction, l'altitude est indiquée par  $J(\theta_0, \theta_1)$  et les pas sont indiqués par le taux d'apprentissage (*Learning rate*).

Pour un point donné, la valeur du gradient (la dérivée) correspond à l'inclinaison de la pente en ce point. La figure suivante illustre la tangente de deux points. Nous pouvons constater que la pente pour le point vert est plus faible que pour le point bleu. Cela indique qu'il y a moins de pas à faire pour atteindre le minimum. Ainsi si le gradient est élevé, c'est que la pente est très pentue, s'il est faible, la pente est faible. S'il est égal à zéro il s'agit d'un minimum. Enfin s'il est négatif cela veut dire que la pente descend (dans ce cas il faut poursuivre dans cette direction) et autrement elle monte (il faut partir dans l'autre direction).



Le point bleu, dans la figure suivante, est situé à la position  $x = -4.5$ . Nous savons que la fonction de coût est :  $y = f(x^2)$ , i.e.  $y = f(-4.5^2) = 20.25$ . La dérivée de  $f$  est :  $f' = 2x$  donc nous avons :  $f'(2 \times -4.5) = -9$ . Le  $-9$  indique qu'il s'agit d'une pente forte et le fait que le nombre soit négatif indique qu'il faut continuer à se déplacer dans cette direction.



Pour résumer le gradient, i.e. la pente de la fonction de coût pour un point donné, représente la direction et le taux de variation de la fonction de coût. Suivre le gradient négatif de la fonction permet donc de la minimiser le plus rapidement possible.

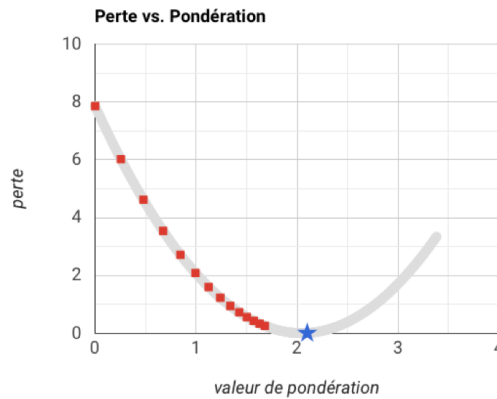
Le pas pour atteindre le minimum s'appelle le taux d'apprentissage (*learning rate*) et est souvent noté  $\eta$ . Plus ce pas est petit plus il faudra du temps pour atteindre le point le plus bas. D'un autre côté un pas trop élevé risque de nous faire dépasser le minimum. Généralement les valeurs de  $\eta$  sont comprises entre 0.0001 et 1.

Le site suivant propose une application pour tester et visualiser les impacts d'un choix d' $\eta$  :

<https://developers.google.com/machine-learning/crash-course/fitter/graph>  
(<https://developers.google.com/machine-learning/crash-course/fitter/graph>).

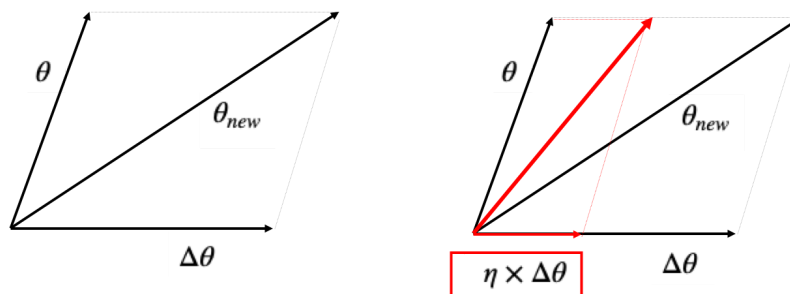
Testez différents taux d'apprentissage et voyez comment ils affectent le nombre d'étapes nécessaires pour atteindre le point minimal de la courbe de coût. Essayez d'effectuer les exercices proposés sous le graph.

|                                   |  |     |
|-----------------------------------|--|-----|
| Définir le taux d'apprentissage : | <input type="range" value="0.2"/>            | 0,2 |
| Exécuter une seule étape :        | <input type="button" value="ÉTAPE"/>         | 13  |
| Réinitialiser le graphe :         | <input type="button" value="RÉINITIALISER"/> |     |



Il est très difficile de choisir la bonne valeur et de nombreuses expériences sont souvent nécessaires. Nous pouvons couvrir plus de domaines avec des pas plus importants (taux d'apprentissage élevé), mais nous courons le risque de dépasser les minima. D'un autre côté, les petits pas (petit taux d'apprentissage) prendront beaucoup de temps pour atteindre le point le plus bas.

Concrètement voici à quoi cela correspond. Si le vecteur des paramètres est  $\theta = [\theta_0, \theta_1]$  et que les changements désirés sont  $\Delta\theta = [\Delta\theta_0, \Delta\theta_1]$  alors la nouvelle position correspond à  $\theta_{new} = \theta_{old} + \eta \times \Delta\theta_{old}$ .



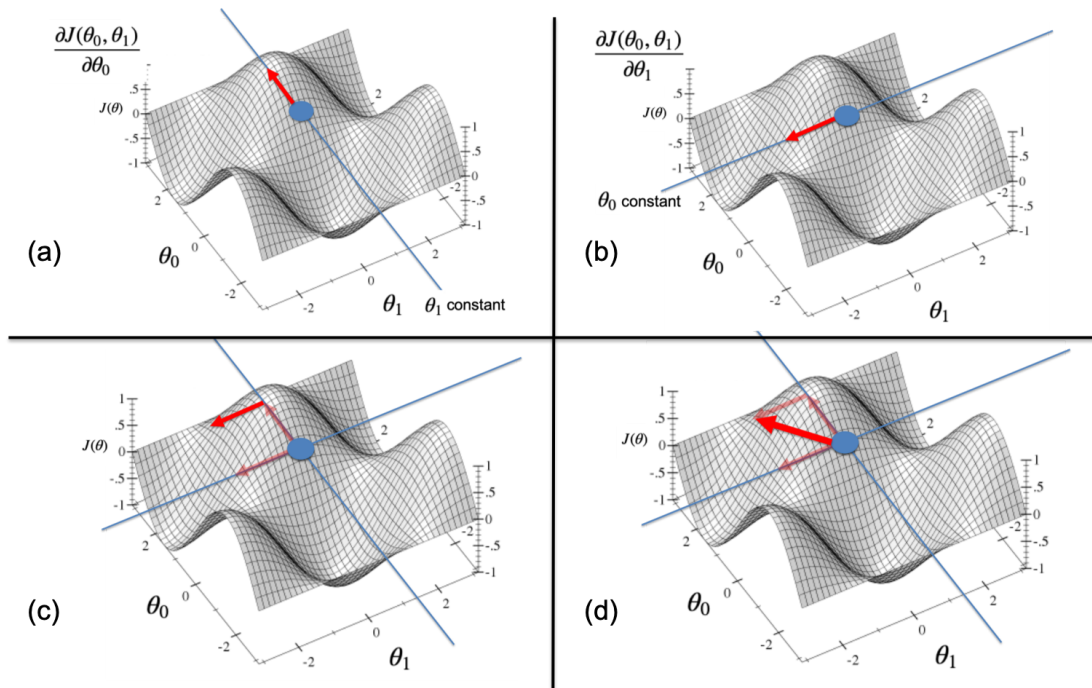
### Petit rappel sur les dérivées partielles

Nous savons que notre fonction de perte est :

$$J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m ((\theta_0 + \theta_1 x^{(i)}) - y^{(i)})^2$$

Le gradient est le vecteur des dérivées partielles selon chaque dimension. Donc pour calculer la descente de gradient il est nécessaire de calculer les dérivées partielles par rapport à  $\theta_0$  et  $\theta_1$  pour savoir dans quelle direction aller.

Une dérivée partielle d'une fonction de plusieurs variables, ici  $J(\theta)$ , est la dérivée de cette fonction selon une variable, les autres étant considérées constantes. Considérons le point positionné sur la figure suivante. Dans (a), nous fixons  $\theta_1$  comme constante et nous pouvons donc calculer la dérivée partielle  $\frac{\partial J(\theta_0, \theta_1)}{\partial \theta_0}$ . Dans notre exemple, le point est positionné en descente le sens du vecteur est donc vers les valeurs négatives. Dans (b), en fixant  $\theta_0$  comme constante, nous pouvons calculer la dérivée partielle  $\frac{\partial J(\theta_0, \theta_1)}{\partial \theta_1}$  et là aussi comme le point est sur une pente, le sens du vecteur est orienté vers des valeurs négatives. Dans (c), il suffit de faire la somme des deux vecteurs précédents et le vecteur résultat est indiqué dans (d).



Ainsi, à partir de :

$$J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m ((\theta_0 + \theta_1 x^{(i)}) - y^{(i)})^2$$

les dérivées partielles par rapport à  $\theta_0$  et  $\theta_1$  :

$$\frac{\partial J(\theta_0, \theta_1)}{\partial \theta_1} = \frac{-2}{m} \sum_{i=1}^m ((\theta_0 + \theta_1 x^{(i)}) - y^{(i)}) x^{(i)}$$

et

$$\frac{\partial J(\theta_0, \theta_1)}{\partial \theta_0} = \frac{-2}{m} \sum_{i=1}^m ((\theta_0 + \theta_1 x^{(i)}) - y^{(i)})$$

L'algorithme de descente de gradient consiste donc à déplacer  $\theta_0$  et  $\theta_1$  jusqu'à atteindre un minimum pour la fonction de coût :

*Initialisation de  $\theta_0$  et  $\theta_1$  avec des nombres aléatoires*

*Répéter jusqu'à convergence*

$$\theta_0 = \theta_0 - \eta \frac{-2}{m} \sum_{i=1}^m ((\theta_0 + \theta_1 x^{(i)}) - y^{(i)})$$

$$\theta_1 = \theta_1 - \eta \frac{-2}{m} \sum_{i=1}^m ((\theta_0 + \theta_1 x^{(i)}) - y^{(i)}) x^{(i)}$$

}

Cela consiste donc à boucler suffisamment longtemps, i.e. spécification du nombre d'epochs par exemple, sur les données pour mettre à jour  $\theta_0$  et  $\theta_1$ .

## Implémentation simple

On redéfinit la fonction predict

In [3]:

```

1  def predict(X, theta0, theta1):
2      return theta0 + theta1*X

```

Définition de la fonction de coût (ici MSE)

In [4]:

```

1  def cost_function(X, y, theta0, theta1):
2      m = len(y)
3      total_error = 0.0
4      for i in range(m):
5          total_error += ((theta0 + theta1*X[i]) - y[i])**2
6      return total_error / float(m)

```

In [5]:

```

1  def update_weights(X, y, theta0, theta1, eta):
2      theta0_deriv = 0
3      theta1_deriv = 0
4      m = len(y)
5
6      for i in range(m):
7          # Calcul des dérivées partielles
8          # -2*((theta0 + theta1*x) - y)
9          theta0_deriv -= -2*((theta0 + theta1*X[i]) - y[i])
10
11         # -2x*((theta0 + theta1*x) - y)
12         theta1_deriv -= -2*X[i] * ((theta0 + theta1*X[i]) - y[i])
13
14         # soustraction ds dérivées et multiplication par eta
15
16         theta0 -= (theta0_deriv / float(m)) * eta
17         theta1 -= (theta1_deriv / float(m)) * eta
18
19     return theta0, theta1
20

```

In [6]:

```

1  def gradient_descent(X, y, theta0, theta1, eta=0.01, epochs=1000):
2
3      #variables pour affichage à la sortie
4      cost_history = []
5      theta_history = np.zeros((epochs,2),dtype=float)
6
7      for i in range(epochs):
8
9          theta0,theta1 = update_weights(X, y, theta0, theta1, eta)
10         # stockage des valeurs pour l'affichage
11         theta_history[i][0]=theta0
12         theta_history[i][1]=theta1
13
14         cost = cost_function(X, y, theta0, theta1)
15         cost_history.append(cost)
16
17     return theta0, theta1, cost_history, theta_history

```

Le code suivant execute l'appel à la fonction

In [7]:

```

1  eta =0.1
2  epochs = 1000
3  X = np.random.rand(100,1)
4  y = 5 + 4*X+np.random.randn(100,1)
5  #initialisation aléatoire de theta0 et de theta1
6  theta0=np.random.randn()
7  theta1=np.random.randn()
8
9
10 theta0,theta1,cost_history,theta_history=gradient_descent(X, y,theta0, theta1
11
12 print ( 'theta0 = %0.2f'%theta0, ' theta1 = %0.2f'%theta1, ' cost(%0.2f) '%cost_hi

```

theta0 = 5.20 theta1 = 3.63 cost(0.83)

In [8]:

```

1  def plot_history (eta,epochs,cost_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      ax.plot(range(epochs),cost_history, 'b.')

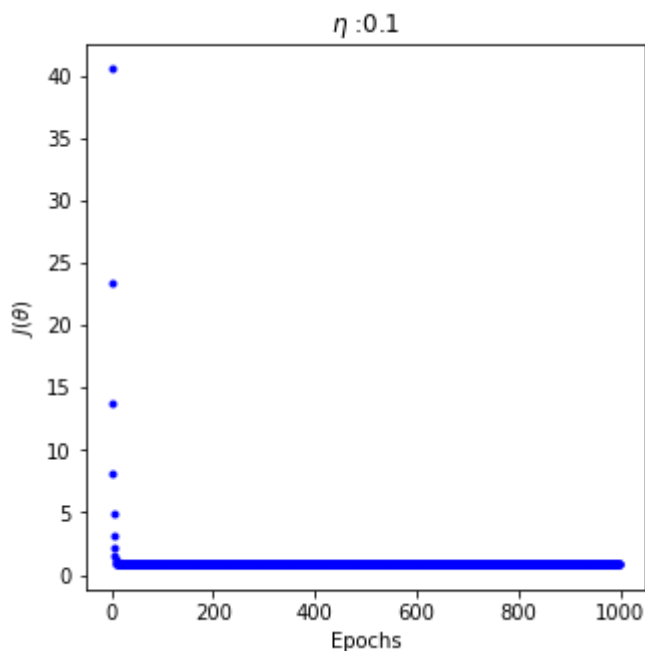
```

In [9]:

```

1  plot_history(eta,epochs,cost_history)

```





In [10]:

```

1  def plot_2D(theta,theta_history, epochs, pas=25):
2      """
3      Affichage en 2D de la cost_function
4      Paramètres :
5      Entrée :
6      -----
7      theta : array theta0 = theta[0], theta1 = theta[1]
8      theta_history array : contient l'historique des thetas lors de l'apprentissage
9      epochs int : epochs utilisés lors de l'apprentissage
10     pas int : indique dans l'historique combien de points doivent être plotés
11     """
12     theta0_reel = theta[0]
13     theta1_reel = theta[1]
14     fig, ax = plt.subplots(figsize=(5,5))
15     # x et y à tracer
16     m=10
17     x = np.linspace(-1,1,m)
18     y = theta0_reel + theta1_reel * x
19
20     #fonction plus simple qui calcule le coût pour les valeurs de l'affichage
21     def cost_func(theta0, theta1):
22         return np.average((y-(theta0 + theta1*x))**2, axis=2)/2
23
24     # Construction de la grille
25     # y_theta0=int(theta0_true+2) pour centrer le dessin idem pour y_theta1
26     # nb_points = nombre de points de la courbe
27     y_theta0=int(theta0_reel+3)
28     y_theta1=int(theta1_reel+3)
29     nb_points=100
30     theta0_grid = np.linspace(-2,y_theta0,nb_points)
31     theta1_grid = np.linspace(-2,y_theta1,nb_points)
32
33     #recuperation de la cost function pour les differentes valeurs de theta0
34     #définie par les theta0_grid et theta1_grid
35
36     J_grid = cost_func(theta0_grid[np.newaxis,:,np.newaxis],
37                        theta1_grid[:,np.newaxis,np.newaxis])
38
39     # Mise en place des courbes de niveau et du contour
40     X, Y = np.meshgrid(theta0_grid, theta1_grid)
41     contours = ax.contour(X, Y, J_grid, 30)
42     ax.clabel(contours)
43     ax.scatter([theta0_reel]*2,[theta1_reel]*2,s=[50,10], color=['b','r'])
44
45     #selection des points à afficher ici tous les pas dans l'historique
46     #on commence par 00 pour le sens des flèches
47     thetas = [np.array((0,0))]
48     N=1
49     for i in range(epochs-1):
50         if i % pas ==0:
51             N+=1
52             this_theta = np.empty((2,))
53             this_theta=theta_history[i]
54             thetas.append(this_theta)
55
56     # fleches
57     for i in range(1,N):
58         ax.annotate('', xy=thetas[i], xytext=thetas[i-1],
59                     arrowprops={'arrowstyle': '->', 'color': 'r', 'lw': 1},

```

```

60         va='center', ha='center')
61
62     ax.scatter(*zip(*thetas), c='b', s=40, lw=0)
63
64     ax.set_xlabel(r'$\theta_0$')
65     ax.set_ylabel(r'$\theta_1$')
66     ax.set_title('Cost function')
67
68     plt.show()

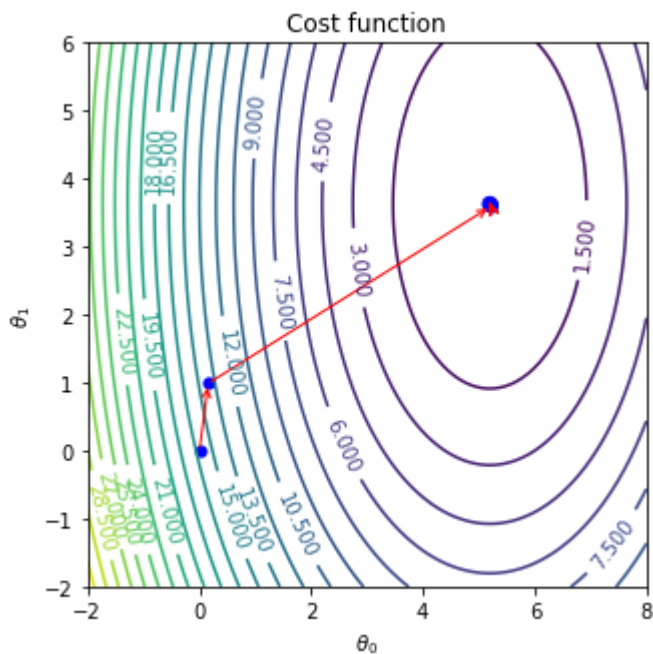
```

In [11]:

```

1  #transformation en tableau pour theta0 et theta1
2  theta=[theta0,theta1]
3  plot_2D(theta,theta_history,epochs,pas=100)

```



In [12]:

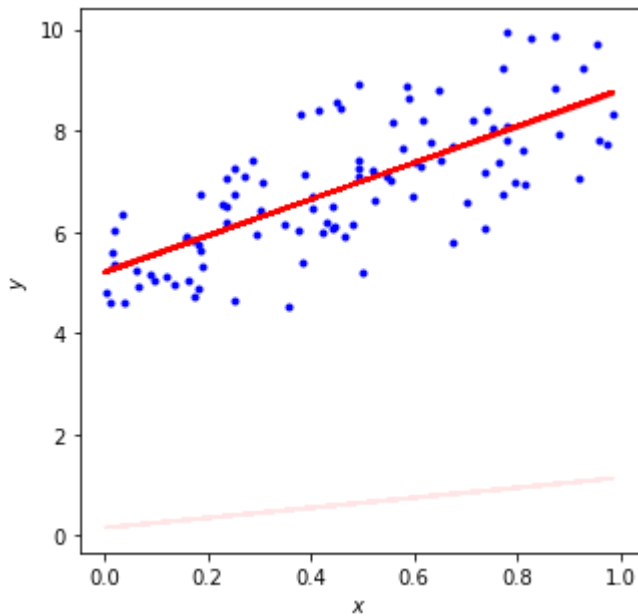
```

1  def plot_lines (eta,epochs,theta_history):
2      tr =0.1
3      fig,ax = plt.subplots(figsize=(5,5))
4      ax.set_xlabel("$x$", fontsize=10)
5      ax.set_ylabel("$y$", fontsize=10)
6      ax.plot(X,y, 'b. ')
7      for i in range (epochs):
8          pred=predict(X, theta_history[i][0], theta_history[i][1])
9          if ((i % 25 == 0) ):
10             _ = ax.plot(X,pred, 'r-',alpha=tr)
11             if tr < 0.8:
12                 tr = tr+0.2
13

```

In [13]:

```
1 plot_lines(eta,epochs,theta_history)
```



In [14]:

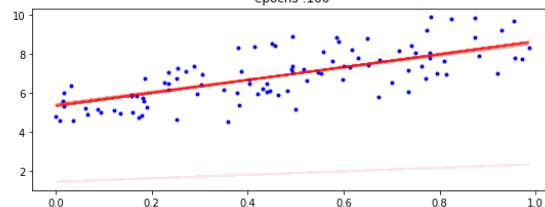
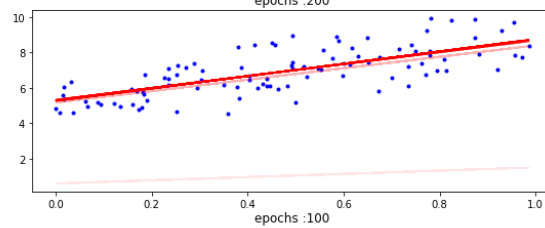
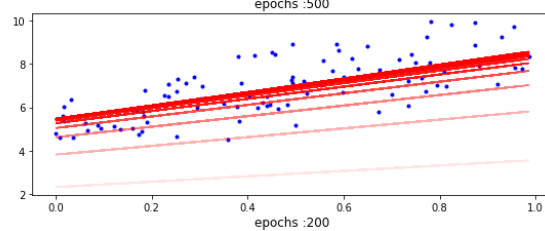
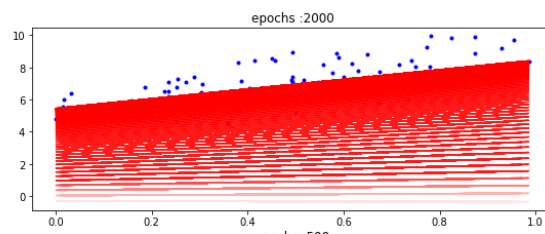
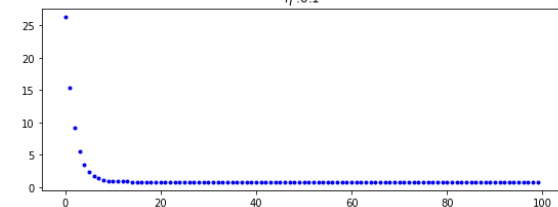
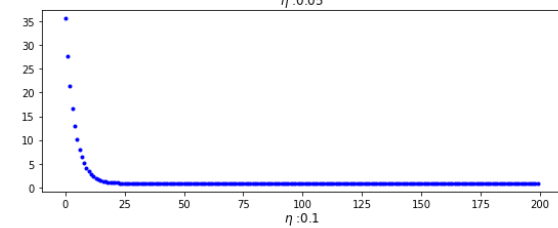
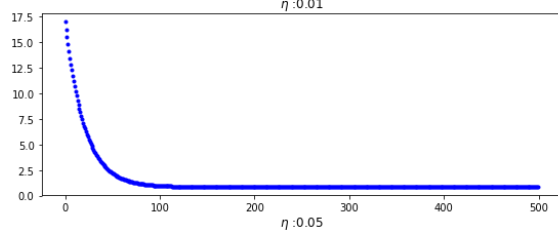
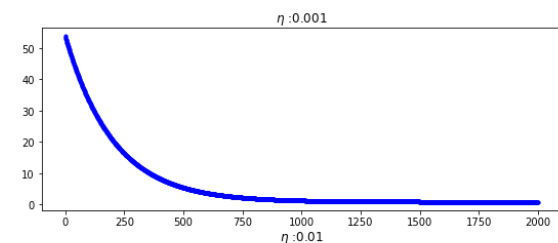
```
1 def plot_all (eta,epochs,theta_history,ax,ax1):
2     tr =0.1
3     #fig,ax = plt.subplots(figsize=(5,5))
4     ax1.plot(X,y,'b.')
5     for i in range (epochs):
6         pred=predict(X, theta_history[i][0], theta_history[i][1])
7         if ((i % 25 == 0) ):
8             ax1.plot(X,pred,'r-',alpha=tr)
9             if tr < 0.8:
10                 tr = tr+0.2
11     ax.plot(range(epochs),cost_history,'b.')
```

In [15]:

```

1  fig = plt.figure(figsize=(20,15))
2  fig.subplots_adjust(hspace=0.2, wspace=0.2)
3  epochs_eta = [(2000,0.001),(500,0.01),(200,0.05),(100,0.1)]
4
5  count =0
6  ▼ for epochs, eta in epochs_eta:
7      count += 1
8      ax = fig.add_subplot(4,2,count)
9      ax.set_title(r"$\eta$ : {}".format(eta))
10     count += 1
11
12     ax1 = fig.add_subplot(4,2,count)
13     ax1.set_title("epochs : {}".format(epochs))
14
15
16     theta0=np.random.randn()
17     theta1=np.random.randn()
18
19     theta0,theta1,cost_history,theta_history=gradient_descent(X, y,theta0, th
20     plot_all(eta,epochs,theta_history,ax,ax1)
21
22

```



## Generalisation et regression linéaire multiple

Dans la section précédente nous avons vu les principes d'une regression linéaire simple. Dans cette section, nous étendons le principe à une régression linéaire multiple (*Multivariate Linear Regression*), i.e. ayant plusieurs paramètres. De plus nous présentons une manière plus efficace d'implémenter la descente de gradient à l'aide de matrice.

Dans une regression linéaire simple on cherche à modéliser la relation entre deux variables quantitatives continues. Comme nous l'avons vu précédemment un modèle de régression linéaire simple est de la forme  $h_\theta(x) = \theta_0 + \theta_1 x$  où :

- $h_\theta(x)$  est le modèle/l'hypothèse pour la variable  $y$  à prédire;
- $x$  est la variable prédictive;
- $\theta_0$  et  $\theta_1$  les paramètres à estimer.

Plus formellement le modèle est plutôt de la forme  $h_\theta(x) = \theta_0 + \theta_1 x + \epsilon$  où  $\epsilon$  correspond au terme d'erreur aléatoire du modèle. Dans notre cas nous le négligeons pour le moment.

Dans une regression linéaire multiple on cherche à modéliser la relation entre plus de deux variables quantitatives continues. Un modèle de régression linéaire multiple est donc de la forme :

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_n x_n$$

Par exemple pour connaître la concentration d'ozone dans l'air ( $y$ ), il faut tenir compte de la température ( $x_1$ ) et de la vitesse du vent ( $x_2$ ).

### Représentation sous la forme de matrice

Ce modèle peut tout à fait être représenté sous la forme de matrices (cette implémentation est beaucoup plus efficace que la précédente) et d'utiliser des produits scalaires entre matrices.

Pour cela il faut au préalable ajouter une colonne de 1 à la matrice des variable prédictives (le nombre de valeurs prédictives est donc incrémenté de 1). Ce 1 permet de pouvoir traiter en même temps que les autres le  $\theta_0$ , i.e.  $h_\theta(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_n x_n$ .

Avec  $m$  le nombre d'enregistrements et  $n$  le nombre de variables prédictives nous avons :

$$h_\theta(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_n x_n \\ = \theta^T X$$

où  $\theta^T$  correspond à la transposée de la matrice  $\theta$ .

$$X = \begin{bmatrix} 1 & x_1^1 & x_2^1 & \dots & x_n^1 \\ 1 & x_1^2 & x_2^2 & \dots & x_n^2 \\ 1 & x_1^3 & x_2^3 & \dots & x_n^3 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^m & x_2^m & \dots & x_n^m \end{bmatrix} \quad (\text{taille : } m \times n + 1)$$

pour les  $\theta$  :

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_2 \\ \theta_3 \\ \vdots \\ \theta_n \end{bmatrix} \text{ (taille : } n + 1 \times 1 \text{)}$$

Et pour  $y$ , nous avons :

$$y = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_m \end{bmatrix} \text{ (taille : } m \times 1 \text{)}$$

La prédiction peut donc, par exemple, facilement se faire par un simple produit scalaire :

$$\begin{array}{c} \xrightarrow{\hspace{10em}} \\ \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_m \end{bmatrix} = \begin{bmatrix} 1 & x_1^1 & x_2^1 & \dots & x_n^1 \\ 1 & x_1^2 & x_2^2 & \dots & x_n^2 \\ 1 & x_1^3 & x_2^3 & \dots & x_n^3 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^m & x_2^m & \dots & x_n^m \end{bmatrix} \cdot \begin{bmatrix} \theta_0 \\ \theta_2 \\ \theta_3 \\ \vdots \\ \theta_n \end{bmatrix} \downarrow \\ \\ \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_m \end{bmatrix} = \begin{bmatrix} 1 & x_1^1 & x_2^1 & \dots & x_n^1 \\ 1 & x_1^2 & x_2^2 & \dots & x_n^2 \\ 1 & x_1^3 & x_2^3 & \dots & x_n^3 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^m & x_2^m & \dots & x_n^m \end{bmatrix} \cdot \begin{bmatrix} \theta_0 \\ \theta_2 \\ \theta_3 \\ \vdots \\ \theta_n \end{bmatrix} \\ \\ \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_m \end{bmatrix} = \begin{bmatrix} \theta_0 & \theta_1 x_1^1 & x_2^1 & \dots & \theta_n x_n^1 \\ \theta_0 & \theta_1 x_1^2 & x_2^2 & \dots & \theta_n x_n^2 \\ \theta_0 & \theta_1 x_1^3 & x_2^3 & \dots & \theta_n x_n^3 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \theta_0 & \theta_1 x_1^m & x_2^m & \dots & \theta_n x_n^m \end{bmatrix} \end{array}$$

La dérivée partielle par rapport à  $\theta_j$  est maintenant :

$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{-2}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

L'algorithme de descente de gradient peut donc se généraliser en :

*Initialisation de  $\theta_j$  avec des nombres aléatoires*

*Répéter jusqu'à convergence* {

$$\theta_j = \theta_j - \eta \frac{-2}{m} \sum_{i=1}^m (h_{\theta} x^{(i)}) - y^{(i)} x_j^{(i)}$$

*en mettant simultanément  $\theta_j$  à jour pour  $j = (1, \dots, n)$*

}

ce qui donne :

$$\begin{aligned}\theta_0 &= \theta_0 - \eta \frac{-2}{m} \sum_{i=1}^m (h_{\theta} x^{(i)} - y^{(i)}) x_0^{(i)} \text{ avec } x_0^{(i)} = 1 \\ \theta_1 &= \theta_1 - \eta \frac{-2}{m} \sum_{i=1}^m (h_{\theta} x^{(i)} - y^{(i)}) x_1^{(i)} \\ \theta_2 &= \theta_2 - \eta \frac{-2}{m} \sum_{i=1}^m (h_{\theta} x^{(i)} - y^{(i)}) x_2^{(i)} \\ &\vdots\end{aligned}$$

ou tout simplement :

*Initialisation de  $\theta_j$  avec des nombres aléatoires*

*Répéter jusqu'à convergence{*

$$\theta = \theta - \eta \frac{-2}{m} (X \cdot \theta - y) X^T$$

*}*

In [16]:

```
1 import numpy as np
2 matrice=np.random.randint(10, size=(2, 3))
3 print ("une matrice\n",matrice)
4 print ("la transposée de la matrice\n",matrice.T)
```

une matrice

```
[[0 4 1]
 [2 1 8]]
```

la transposée de la matrice

```
[[0 2]
 [4 1]
 [1 8]]
```

## Implémentation

Nous définissons d'abord la fonction predict

In [17]:

```
1 def predict(X, theta):
2     """
3     Prediction de la valeur de X
4     Paramètres :
5     -----
6     Input :
7     X array : données d'apprentissage (variables prédictives)
8     theta array : theta[0], theta1 = theta[1]
9     -----
10    Output :
11    array correspondant au produit matricielle de X avec les thetas
12    """
13    return X.dot(theta)
```

La fonction de coût se fait simplement à l'aide d'un produit scalaire (fonction dot en python)

In [18]:

```
1  def cost_function(X,y,theta):
2      """
3      Calcul de la fonction de coût
4      Paramètres :
5      -----
6      Input :
7      X array : données d'apprentissage (variables prédictives)
8      y array : données d'apprentissage (valeur à prédire)
9      theta array : theta[0], theta1 = theta[1]
10     -----
11     Output :
12     MSE
13     """
14
15     m=len(y)
16     prediction=X.dot(theta)
17     square_err=(prediction - y)**2
18     return 1/(2*m) * np.sum(square_err)
```

La descente de gradient (ne pas oublier la transposé de X)



In [19]:

```

1  def gradient_descent(X,y,theta,eta=0.01, epochs=1000):
2      """
3      Descente de gradient simple (vanilla descent gradient)
4      Paramètres :
5      -----
6      Input :
7      X array : données d'apprentissage (variables prédictives)
8      y array : données d'apprentissage (valeur à prédire)
9      theta array : theta[0], theta1 = theta[1]
10     eta float : le learning rate
11     epochs int : nombre de répétitions
12     -----
13     Output :
14     theta array : theta[0], theta1 = theta[1]
15     cost_history array : liste du coût de theta à chaque itération
16     theta_history array : liste des thetas à chaque itération
17     """
18
19     m=len(y) # attention prendre y et non x car il y a un entier en plus
20     cost_history=[]
21     theta_history = np.zeros((epochs,2),dtype=float)
22     #Attention theta_history est pour afficher en 2D donc vecteur à 2 dimensi
23
24     for i in range(epochs):
25
26         #Le code ci-dessous peut s'écrire en une ligne
27         #theta = theta -(1/(2*m))*eta*( X.T.dot((prediction - y)))
28         prediction = predict(X,theta)
29         error = X.T.dot((prediction - y))
30         descent=eta * (1/(2*m)) * error
31         theta-=descent
32
33         #pour affichage
34         theta_history[i][0]=theta[0]
35         theta_history[i][1]=theta[1]
36         cost_history.append(cost_function(X,y,theta))
37
38     return theta, cost_history, theta_history

```

In [20]:

```

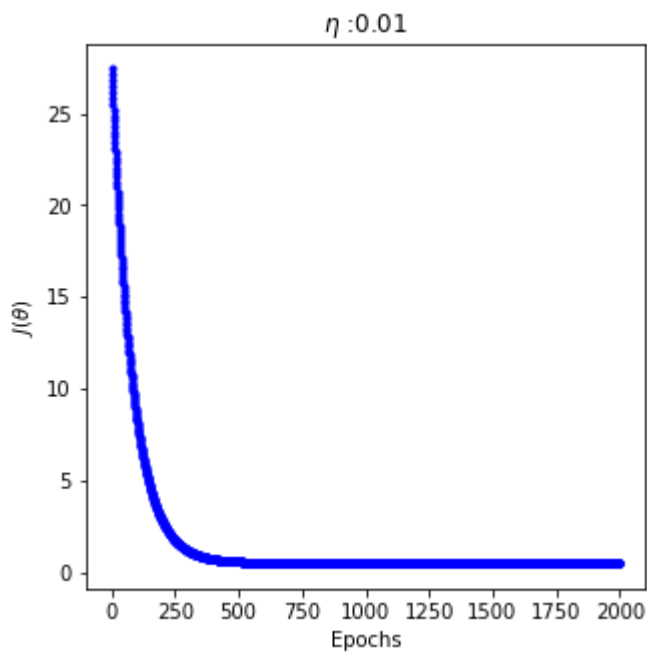
1  eta =0.01
2  epochs = 2000
3  #initialisation aléatoire de theta0 et de theta1
4  # dans cette exemple on considère une fonction du type
5  # theta0 + theta1x_1 pour l'affichage en 2D elle est bien sûr généralisable
6  theta = np.random.randn(2,1)
7  X = np.random.rand(100,1)
8  y = 5 + 4*X+np.random.randn(100,1)
9
10 #X_0 ajoute un 1 au début de X pour avoir : theta0x_0 + theta1x_1 ...
11 X_0 = np.c_[np.ones((X.shape[0])), X]
12
13
14 #appel de la fonction attention c'est avec X_0
15 theta,cost_history,theta_history = gradient_descent(X_0,y,theta,eta,epochs)
16 print ('theta0 = %0.2f'%theta[0], ' theta1 = %0.2f'%theta[1], ' cost(%0.2f)'%cos
17

```

theta0 = 5.31   theta1 = 3.76   cost(0.45)

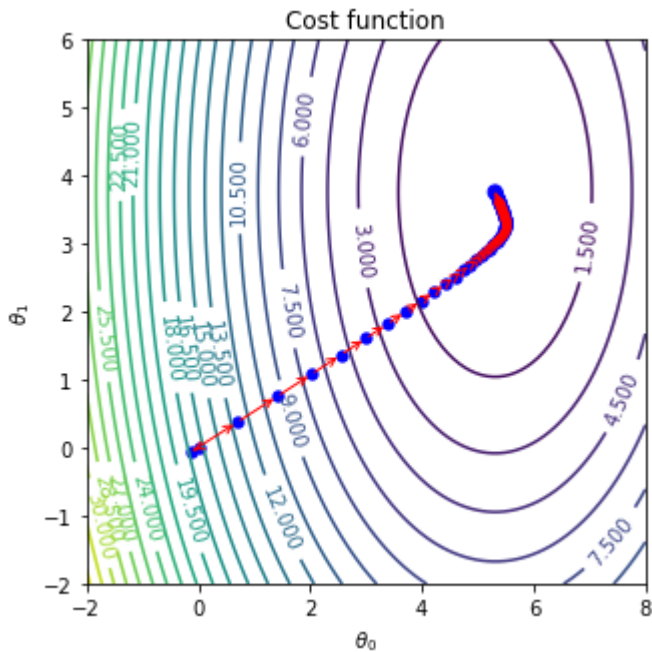
In [21]:

```
1 plot_history(eta,epochs,cost_history)
```



In [22]:

```
1 plot_2D(theta,theta_history,epochs,pas=25)
```



## Optimisation : différentes stratégies

L'algorithme de descente de gradient que nous avons vu précédemment est également appelée descente de gradient par lots (*Batch Gradient Descent*) et est le plus simple : il calcule l'erreur pour chaque exemple dans l'ensemble d'apprentissage. Après avoir évalué tous les enregistrements de l'ensemble d'apprentissage (1 epoch), il met à jour les paramètres du modèle.

Il présente l'avantage d'être assez efficace en termes de calcul et de produire un gradient d'erreur stable et une convergence stable. L'un des inconvénients est que ce gradient d'erreur stable peut parfois conduire à un état de convergence qui n'est pas optimal. L'autre inconvénient (important) est aussi qu'il nécessite que l'ensemble du jeu de données d'apprentissage réside en mémoire et ce n'est pas possible pour des jeux de données volumineux.

Pour cela différentes stratégies ont été proposées.

### Stochastic gradient descent

La descente de gradient stochastique (*Stochastic Gradient Descent*) met à jour les paramètres en fonction du gradient de l'erreur par rapport à un seul exemple d'apprentissage (`batch_size=1`) tiré aléatoirement (d'où le terme stochastique) même s'il est tout à fait possible de considérer un sous ensemble d'exemples. Ceci est différent de Batch Gradient Descent, qui met à jour les paramètres une fois que tous les exemples d'apprentissage ont été évalués.

Il ne consomme pas beaucoup de mémoire et est souvent plus rapide car un seul échantillon est traité à la fois. Ainsi pour des grands ensembles de données, la convergence peut être plus rapide car les mises à jour des paramètres sont plus fréquentes. L'avantage de ces mises à jour fréquentes est que les étapes effectuées vers les minima de la fonction de perte ont des oscillations qui peuvent aider à sortir des minimums locaux.

En fonction des données, elle est généralement plus rapide que la descente de gradient par lots. Cependant il s'avère qu'à cause des mises à jour fréquentes la convergence prenne plus de temps (descente dans des directions différentes). Les mises à jour fréquentes sont coûteuses en termes de calcul car elles utilisent toutes les ressources pour traiter un échantillon d'apprentissage à la fois. Enfin, quand il y a un seul exemple, nous perdons tous les avantages de l'utilisation des matrices.

### Implémentation du Stochastique Gradient Descent

Dans la boucle principale (epoch), lors du parcours du jeu d'apprentissage, un nombre aléatoire d'exemples sont choisis et le calcul se fait sur eux.

In [23]:

```

1  def stocashtic_gradient_descent(X,y,theta,learning_rate=0.01,epochs=1000):
2      """
3      Descente de gradient stochastique
4      Paramètres :
5      -----
6      Input :
7      X array : données d'apprentissage (variables prédictives) avec un 1 ajout
8      y array : données d'apprentissage (valeur à prédire)
9      theta array : theta[0], theta1 = theta[1]
10     eta float : le learning rate
11     epochs int : nombre de répétitions
12     -----
13     Output :
14     theta array : theta[0], theta1 = theta[1]
15     cost_history array : liste du coût de theta à chaque itération
16     theta_history array : liste des thetas à chaque itération
17     """
18
19     m=len(y)
20     cost_history=[]
21     theta_history = np.zeros((epochs,2),dtype=float)
22     for i in range(epochs):
23         cost = 0.0
24         for j in range(m):
25             #tirage d'un nombre aléatoire entre 0 et m
26             rand_ind = np.random.randint(0,m)
27             #recuperation du X et y correspondants au nombre aléatoire
28             X_j = X[rand_ind,:].reshape(1,X.shape[1])
29             y_j = y[rand_ind].reshape(1,1)
30             prediction = predict(X_j,theta)
31             error = X_j.T.dot((prediction - y_j))
32             descent=eta * (1/(2*m)) * error
33             theta-=descent
34             cost+=cost_function(X_j,y_j,theta)
35
36             theta_history[i][0]=theta[0]
37             theta_history[i][1]=theta[1]
38             cost_history.append(cost)
39     return theta, cost_history, theta_history

```

In [24]:

```

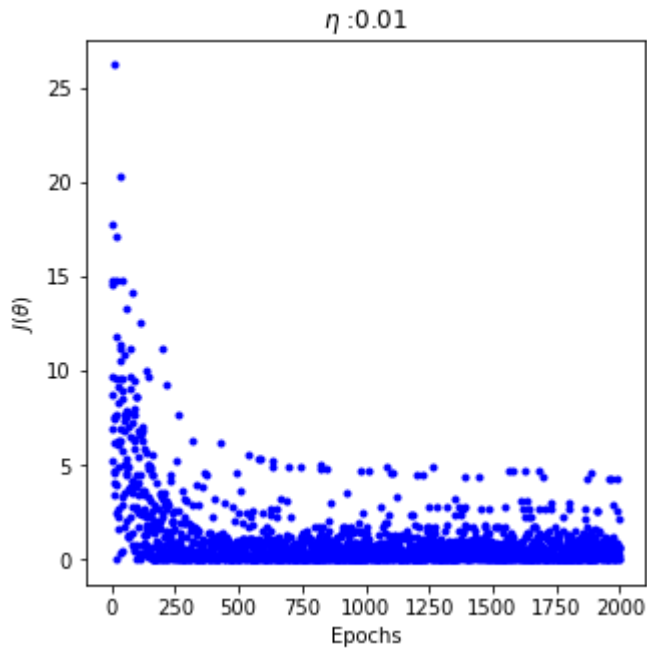
1  eta =0.01
2  epochs = 2000
3  #initialisation aléatoire de theta0 et de theta1
4  theta = np.random.randn(2,1)
5  X = np.random.rand(100,1)
6  y = 5 + 4*X+np.random.randn(100,1)
7
8  #X_0 ajoute un 1 au début de X pour avoir : theta0x_0 + theta1x_1 ...
9  X_0 = np.c_[np.ones((X.shape[0])), X]
10
11  theta,histroy_stoc,theta_history = stocashtic_gradient_descent(X_0,y,theta,et
12  print ('theta0 = %0.2f'%theta[0], ' theta1 = %0.2f'%theta[1], ' cost(%0.2f)'%his

```

theta0 = 5.38   theta1 = 2.96   cost(0.63)

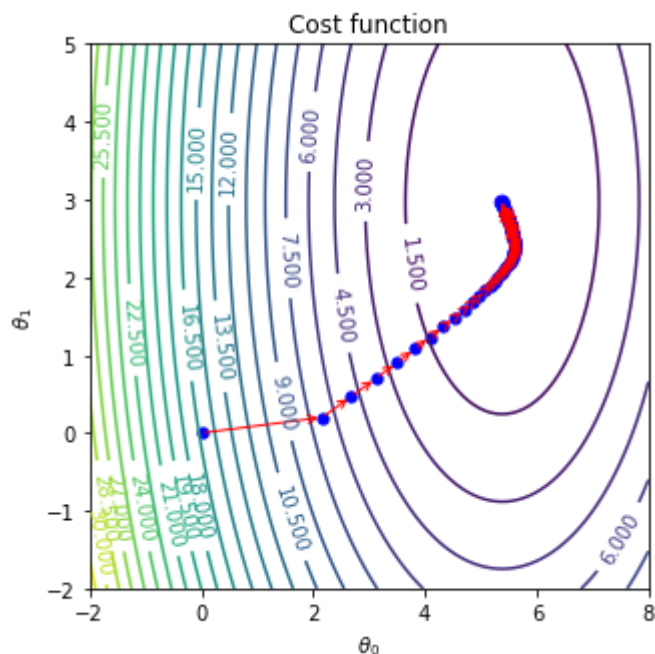
In [25]:

```
1 plot_history(eta,epochs,history_stoc)
```



In [26]:

```
1 plot_2D(theta,theta_history,epochs,pas=25)
```



### Mini-batch gradient descent

L'algorithme de descente de gradient par mini lots (*mini-batch gradient descent*) est souvent privilégié, car il associe une descente en gradient stochastique et une descente en gradient par lots. Il sépare simplement le jeu d'apprentissage en petits lots et effectue une mise à jour pour chacun de ces lots. Cela crée donc un équilibre entre l'efficacité de la descente de gradient par lots et la robustesse de la descente de gradient stochastique. Il est très utilisé en apprentissage profond (deep learning).

### Implémentation du Mini-batch gradient descent

Le jeu d'apprentissage est découpé en petit batch en tirant des valeurs d'exemple au hasard. Ensuite chaque batch est examiné par la descente de gradient.

In [27]:

```
1  ▼ def next_batch(X, y, batchsize):
2      """
3      Création des batches
4      Paramètres :
5      -----
6      Input :
7      X array : données d'apprentissage (variables prédictives) avec un 1 ajout
8      y array : données d'apprentissage (valeur à prédire)
9      batchsize int : taille du batch
10     """
11     m=len(y)
12     # permutation aléatoire de X et y pour faire des batchs avec des valeurs
13     indices = np.random.permutation(m)
14     X = X[indices]
15     y = y[indices]
16     ▼ for i in np.arange(0, X.shape[0], batchsize):
17         # creation des batchs de taille batchsize
18         yield (X[i:i + batchsize], y[i:i + batchsize])
```

In [28]:

```

1  def minibatch_gradient_descent(X,y,theta,learning_rate=0.01,epochs=10,batchsi
2  """
3  Descente de gradient par mini batch
4  Paramètres :
5  -----
6  Input :
7  X array : données d'apprentissage (variables prédictives) avec un 1 ajout
8  y array : données d'apprentissage (valeur à prédire)
9  theta array : theta[0], theta1 = theta[1]
10 eta float : le learning rate
11 epochs int : nombre de répétitions
12 batchsize int : taille du batch
13 -----
14 Output :
15 theta array : theta[0], theta1 = theta[1]
16 cost_history array : liste du coût de theta à chaque itération
17 theta_history array : liste des thetas à chaque itération
18 """
19 m=len(y)
20 cost_history=[]
21 theta_history = np.zeros((epochs,2),dtype=float)
22 for i in range(epochs):
23     cost_batch = []
24     for (batchX, batchY) in next_batch(X, y, batchsize):
25         # Extraction et traitement d'un batch à la fois
26         prediction = predict(batchX,theta)
27         error = batchX.T.dot((prediction - batchY))
28         descent=eta * (1/(2*m)) * error
29         theta-=descent
30         cost_batch.append(cost_function(batchX,batchY,theta))
31         theta_history[i][0]=theta[0]
32         theta_history[i][1]=theta[1]
33
34     cost_history.append(np.average(cost_batch))
35
36 return theta, cost_history, theta_history
37

```

In [29]:

```

1  eta =0.1
2  epochs = 2000
3  #initialisation aléatoire de theta0 et de theta1
4  theta = np.random.randn(2,1)
5  X = np.random.rand(100,1)
6  y = 5 + 4*X+np.random.randn(100,1)
7
8  #X_0 ajoute un 1 au début de X pour avoir : theta0x_0 + theta1x_1 ...
9  X_0 = np.c_[np.ones((X.shape[0])), X]
10
11 theta,history_mini,theta_history = minibatch_gradient_descent(X_0,y,theta,eta
12 print ('theta0 = %0.2f'%theta[0], ' theta1 = %.2f'%theta[1], ' cost(%0.2f)'%his

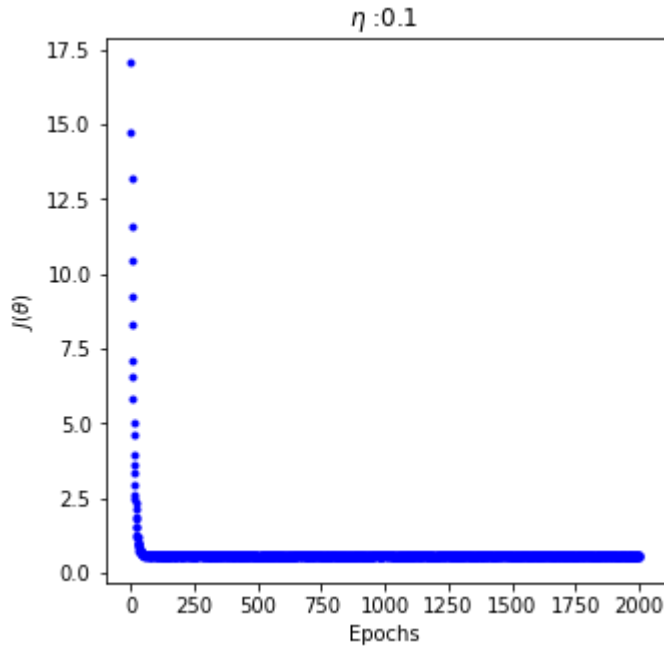
```

theta0 = 4.89    theta1 = 3.99    cost(0.59)



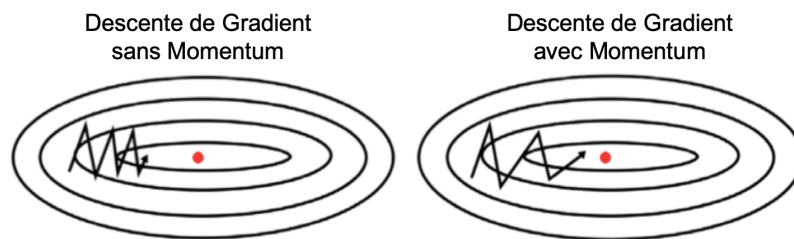
In [30]:

```
1 plot_history(eta, epochs, history_mini)
```



## Momentum

Les approches précédentes ont du mal dans les zones où la surface s'incline beaucoup plus fortement dans une dimension que dans une autre. Dans ces scénarios, la descente de gradient oscille sur les pentes et ne progresse que doucement vers le minima local. Momentum est une méthode qui permet d'accélérer la descente de gradient dans la direction voulue et d'atténuer les oscillations (C.f Figure).



Le principe est le suivant. Considérons une balle qui descend le long d'une colline. Elle accumule de l'élan (*Momentum*) dans la descente et devient de plus en plus rapide jusqu'à ce qu'elle atteigne sa vitesse limite s'il existe une résistance de l'air. Dans le cas de la mise à jour des paramètres, la quantité de mouvement augmente pour les dimensions dont les gradients pointent dans la même direction et réduit les mises à jour pour les dimensions dont les gradients changent de direction. Par conséquent la convergence est plus rapide car le nombre d'oscillation diminue.

Momentum permet de modifier la descente de gradient en introduisant deux nouveaux paramètres : la vitesse  $v$  que nous essayons d'optimiser et le frottement ( $\beta$ ).

Formellement il est défini, pour une fonction de coût  $J$  quelconque, de la manière suivante :

$$\begin{aligned} v &= \beta v - \eta \Delta J \\ \theta &= \theta - v \end{aligned}$$

Dans le cas des fonctions linéaires précédentes nous avons :

$$v = \beta v - \eta \times \frac{-2}{m} \times (X \cdot \theta - y) X^T$$

$$\theta = \theta - v$$

L'avantage de momentum est qu'il ne change pas beaucoup l'algorithme de descente de gradient mais qu'il accélère considérablement l'apprentissage. Généralement  $\beta$  est fixé à 0.9.

### Implémentation du momentum

L'approche par mini-batch étant la plus efficace, il suffit de remplacer la partie mise à jour des gradients par la formule précédente.

In [31]:

```

1  def next_batch(X, y, batchsize):
2      """
3      Création des batches
4      Paramètres :
5      -----
6      Input :
7      X array : données d'apprentissage (variables prédictives) avec un 1 ajout
8      y array : données d'apprentissage (valeur à prédire)
9      batchsize int : taille du batch
10     """
11     m=len(y)
12     # permutation aléatoire de X et y pour faire des batchs avec des valeurs
13     indices = np.random.permutation(m)
14     X = X[indices]
15     y = y[indices]
16     for i in np.arange(0, X.shape[0], batchsize):
17         # creation des batchs de taille batchsize
18         yield (X[i:i + batchsize], y[i:i + batchsize])

```

In [32]:

```

1  def minibatch_gradient_descent_Momentum(X,y,v,beta,theta,learning_rate=0.01,e
2  """
3  Descente de gradient par mini batch avec momentum
4  Paramètres :
5  -----
6  Input :
7  X array : données d'apprentissage (variables prédictives) avec un 1 ajout
8  y array : données d'apprentissage (valeur à prédire)
9  theta array : theta[0], theta1 = theta[1]
10 eta float : le learning rate
11 epochs int : nombre de répétitions
12 batchsize int : taille du batch
13 beta : momentum
14 -----
15 Output :
16 theta array : theta[0], theta1 = theta[1]
17 cost_history array : liste du coût de theta à chaque itération
18 theta_history array : liste des thetas à chaque itération
19 """
20 m=len(y)
21 cost_history=[]
22 v_history=np.zeros((epochs,2),dtype=float)
23 theta_history = np.zeros((epochs,2),dtype=float)
24
25 #beta=0.9
26 for i in range(epochs):
27     cost_batch = []
28     for (batchX, batchY) in next_batch(X, y, batchsize):
29         # Extraction et traitement d'un batch à la fois
30         prediction = predict(batchX,theta)
31         error = batchX.T.dot((prediction - batchY))
32         # calcul du momentum
33         grad=(1/(2*m)) * error
34         v = beta * v + eta*grad
35         theta-= v
36         cost_batch.append(cost_function(batchX,batchY,theta))
37         theta_history[i][0]=theta[0]
38         theta_history[i][1]=theta[1]
39
40     cost_history.append(np.average(cost_batch))
41
42 return theta, cost_history, theta_history
43
44

```

In [33]:

```

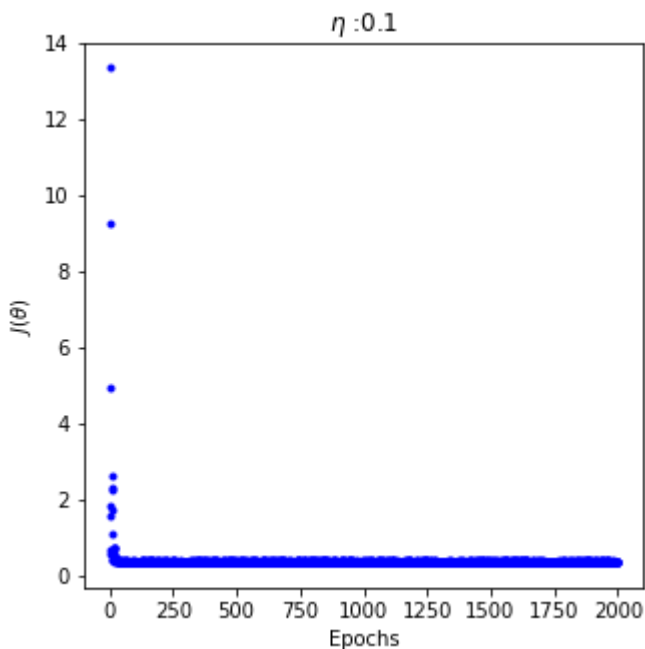
1  eta =0.1
2  epochs = 2000
3  beta = 0.9 # frottement
4  #initialisation aléatoire de theta0 et de theta1
5  theta = np.random.randn(2,1)
6  X = np.random.rand(100,1)
7  y = 5 + 4*X+np.random.randn(100,1)
8  # ajout du vecteur de vitesse
9  v = np.zeros(shape=(2,1))
10 #X_0 ajoute un 1 au début de X pour avoir : theta0x_0 + theta1x_1 ...
11 X_0 = np.c_[np.ones((X.shape[0])), X]
12
13 theta,history_minimomentum,theta_history = minibatch_gradient_descent_Momentu
14 print ('theta0 = %0.2f'%theta[0], ' theta1 = %0.2f'%theta[1], ' cost(%0.2f)'%his

```

theta0 = 4.94 theta1 = 4.08 cost(0.59)

In [34]:

```
1 plot_history(eta,epochs,history_minimomentum)
```



Il existe d'autres types d'optimisation :

- *Nesterov Momentum* assez proche de Momentum mais qui, au lieu de calculer le gradient à la position actuelle, le calcul à une nouvelle position approximé.
- *Adagrad* s'intéresse au learning rate. Précédemment  $\eta$  a toujours été considéré constant. L'objectif d'Adagrad est de s'adapter pour pouvoir accélérer ou ralentir.
- *RMSprop* est une extension d'Adagrad pour corriger le fait que ce dernier considère une fonction d'accumulation du gradient qui est croissante monotone. Ceci peut cependant poser des problèmes car le taux d'apprentissage diminue également de façon monotone et l'apprentissage peut s'arrêter car  $\eta$  devient trop petit.
- *ADAM* cumule les avantages de Momentum et de RMSprop.

Toutes ces optimisations sont disponibles dans la plupart des approches de deep learning et sont utilisées lors des étapes de backpropagation.

## Un exemple d'utilisation : la regression logistique

La régression logistique (*Logistic regression*) ou modèle logit constitue un cas particulier de modèle linéaire généralisé et est particulièrement utilisé comme approche supervisée de classification. Par contre, contrairement au modèle linéaire, la variable prédite ne peut prendre qu'un nombre limité de valeurs.

Dans le cas de la régression logistique binaire, étant donné un ensemble d'entrées  $X$ , l'objectif est de les classer dans l'une des deux catégories possible : 0 ou 1. En d'autres termes, elle modélise la probabilité que chaque entrée appartienne à une catégorie particulière.

Dans notre cas, nous savons que le modèle de régression linéaire est de la forme :

$$h_{\theta}(x) = \theta^T x$$

Pour pouvoir effectuer la classification, nous devons appliquer une fonction :

$$h_{\theta}(x) = \sigma(\theta^T x)$$

Dans le cas binaire, l'une des fonction d'activation couramment utilisée est la fonction sigmoïde

$$\sigma(t) = \frac{1}{1 + e^{-t}}$$

L'hypothèse pour la regression logistique devient donc :

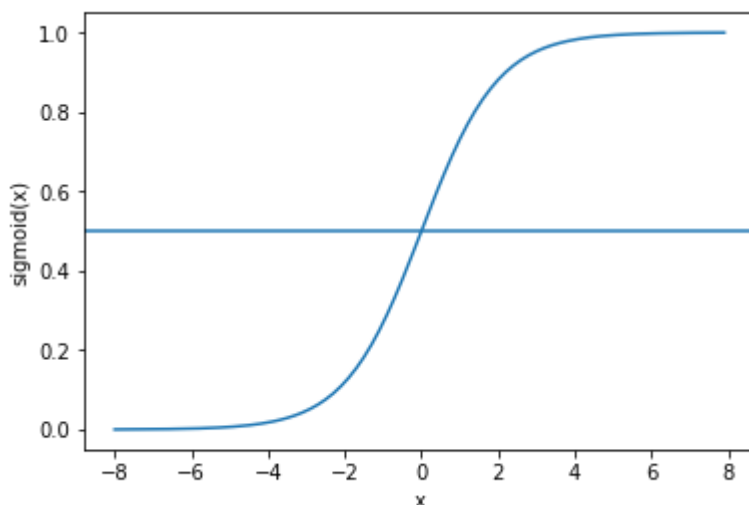
$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

In [35]:

```
1 def sigmoid(z):  
2     return 1 / (1 + np.exp(-z))
```

In [36]:

```
1 x = np.arange(-8, 8, 0.1)  
2 f = sigmoid(x)  
3 plt.plot(x, f)  
4 plt.xlabel('x')  
5 plt.axhline(.5)  
6 plt.ylabel('sigmoid(x)')  
7 plt.show()
```



Nous remarquons sur la courbe que  $\text{sigmoid}(x) > 0.5$  quand  $x > 0$  et  $\text{sigmoid}(x) < 0.5$  quand  $x < 0$  :

$$h_{\theta}(x) = \begin{cases} > 0.5, & \text{si } \theta^T x > 0 \\ < 0.5, & \text{si } \theta^T x < 0 \end{cases}$$

En d'autres termes, si la somme pondérée des entrées est supérieure à zéro, la classe prédite est 1 et inversement si elle est inférieure à 0, la classe prédite est 0.

La fonction de coût généralement associée est la suivante :

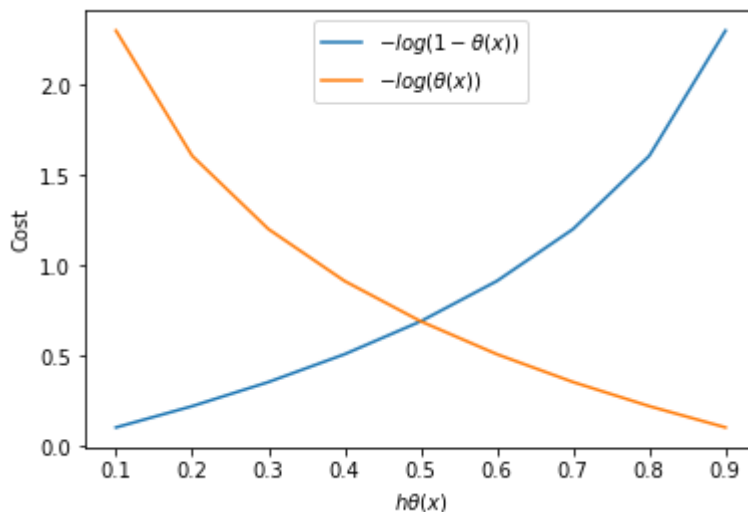
$$\text{Cost} = \begin{cases} -\log(h_{\theta}(x)), & \text{si } y = 1 \\ -\log(1 - h_{\theta}(x)), & \text{si } y = 0 \end{cases}$$

L'intuition derrière la fonction de coût est la suivante :

Si la classe réelle est 1 et que le modèle prédit 0, il faut fortement le pénaliser et vice-versa. Comme vous pouvez le constater sur la figure ci-dessous, pour le graphique  $-\log(h_{\theta}(x))$ , lorsque  $h_{\theta}(x)$  approche de 1, le coût est égal à 0 et lorsque  $h_{\theta}(x)$  est proche de 0, le coût est égal à l'infini (le modèle est fortement pénalisé). De même pour le graphique  $-\log(1 - h_{\theta}(x))$  lorsque la valeur réelle est 0 et que le modèle prédit 0, le coût est 0 et le coût devient l'infini lorsque  $h_{\theta}(x)$  se rapproche de 1.

In [37]:

```
1 x = np.arange(0.1, 1, 0.1)
2 f = -np.log(x)
3 g = -np.log(1-x)
4 plt.plot(x, g, label=r'$-\log(1-\theta(x))$')
5 plt.plot(x, f, label=r'$-\log(\theta(x))$')
6 plt.xlabel(r'$h_{\theta}(x)$')
7 #plt.xlabel(r'$h_{\theta}(x)$')
8 plt.ylabel('Cost')
9 plt.legend()
10 plt.show()
```



Nous pouvons donc combiner les deux équations précédentes :

$$\text{Cost}(h_{\theta}(x), y) = -y \log(h_{\theta}(x)) - (1 - y) \log(1 - h_{\theta}(x))$$

La perte pour tous les exemples d'apprentissage,  $m$ , est comme précédemment calculé en prenant la moyenne de tous les coûts de tous les échantillons d'apprentissage :

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m [-y^{(i)} \log(h_{\theta}(x^{(i)})) - (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)}))]$$

ou

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)}))]$$

La dérivée partielle par rapport à  $\theta_j$  est :

$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y) x_j^{(i)}$$

La descente de gradient consiste donc à répéter, comme précédemment, jusqu'à convergence :

$$\theta_j = \theta_j - \eta \frac{1}{m} X^T (\sigma(X \cdot \theta) - y)$$

In [38]:

```
1  def cost_function(h, y):
2      return (-y * np.log(h) - (1 - y) * np.log(1 - h)).mean()
```

Sélection du jeu de données Iris et récupération des 100 premiers enregistrements (*Serosa* et *Versicolor*). Nous nous intéressons uniquement aux deux premiers attributs (*longueur des sépales* et *largeur des sépales*).

In [39]:

```
1  from sklearn import datasets
2  data = datasets.load_iris()
3
4  #récupération des 100 premiers (ils sont ordonnés)
5  X = data.data[:100, :2]
6  y = data.target[:100]
7
```

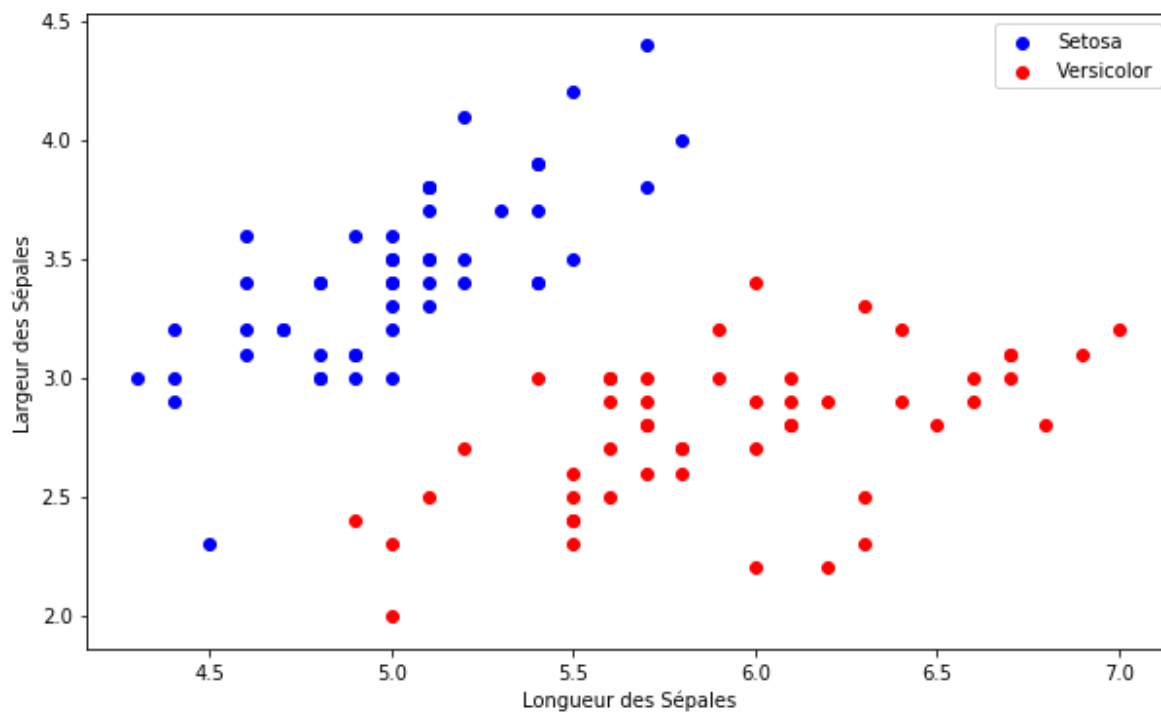
Affichage du jeu de données

In [40]:

```
1 plt.figure(figsize=(10, 6))
2 setosa = plt.scatter(X[:50,0], X[:50,1], c='b')
3 versicolor = plt.scatter(X[50:,0], X[50:,1], c='r')
4 plt.xlabel("Longueur des Sépales")
5 plt.ylabel("Largeur des Sépales")
6 plt.legend((setosa, versicolor), ("Setosa", "Versicolor"))
```

Out[40]:

<matplotlib.legend.Legend at 0x11bc069b0>





In [41]:

```

1  def gradient_descent(X, y, theta, eta=0.01, epochs=1000):
2      m=len(y)
3      cost_history=[]
4      for i in range(epochs):
5          z = np.dot(X, theta)
6          h = sigmoid(z)
7          gradient = np.dot(X.T, (h - y)) / m
8          theta -= eta * gradient
9
10         cost_history.append(cost_function(h, y))
11
12     return theta, cost_history

```

In [42]:

```

1  eta = 0.1
2  epochs = 30000
3  #Initialisation de theta (ici à zéro)
4  theta = np.zeros(X.shape[1])
5
6  theta, cost_history = gradient_descent(X, y, theta, eta, epochs)
7
8  print ("Les paramètres sont : ", theta)

```

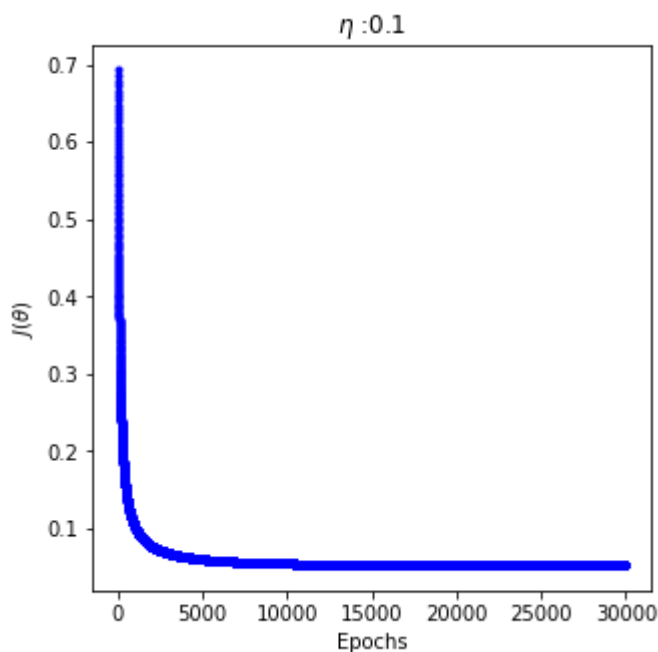
Les paramètres sont : [ 6.66602804 -11.67484425]

In [43]:

```

1  plot_history(eta, epochs, cost_history)

```



Définissons, à présent, la fonction de prediction. Par rapport ce que nous avons vu précédemment, la valeur de seuil pour sigmoide est de 0.5.

In [44]:

```
1 def predicts(X, theta, threshold=0.5):
2     return sigmoid(np.dot(X, theta)) >= threshold
```

In [45]:

```
1 y_predits = predicts(X,theta)
2 print (theta)
3 preds = predicts(X,theta)
4 (preds == y).mean()
5 print ("Nombre de bons prédictions", np.sum(y == y_predits), ' sur ', y.size)
```

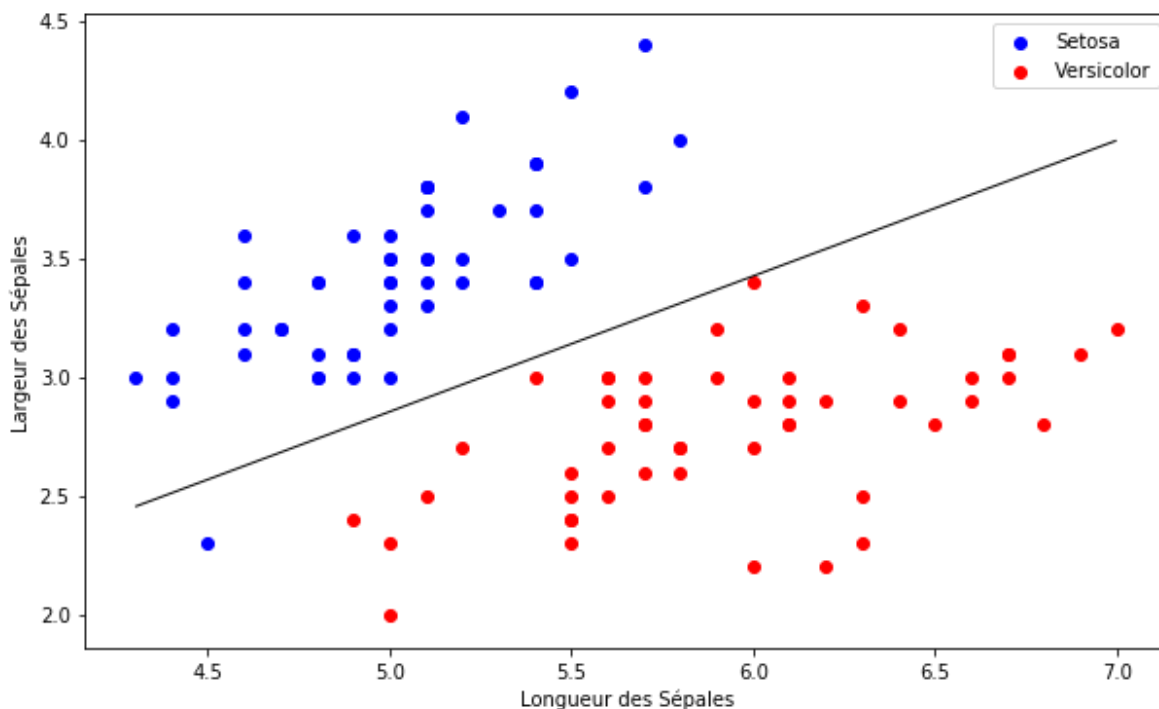
[ 6.66602804 -11.67484425]

Nombre de bons prédictions 99 sur 100

Affichage de la zone de décision

In [46]:

```
1 plt.figure(figsize=(10, 6))
2 setosa = plt.scatter(X[:50,0], X[:50,1], c='b')
3 versicolor = plt.scatter(X[50:,0], X[50:,1], c='r')
4 plt.xlabel("Longueur des Sépales")
5 plt.ylabel("Largeur des Sépales")
6 plt.legend((setosa, versicolor), ("Setosa", "Versicolor"))
7
8
9 x1_min, x1_max = X[:,0].min(), X[:,0].max(),
10 x2_min, x2_max = X[:,1].min(), X[:,1].max(),
11 xx1, xx2 = np.meshgrid(np.linspace(x1_min, x1_max), np.linspace(x2_min, x2_max))
12 grid = np.c_[xx1.ravel(), xx2.ravel()]
13 probs = sigmoid(np.dot(grid, theta)).reshape(xx1.shape)
14 plt.contour(xx1, xx2, probs, [0.5], linewidths=1, colors='black');
```



Essayons à présent le modèle de Regression Logistique de Sickit Learn.

In [47]:

```

1 from sklearn import linear_model
2 clf = linear_model.LogisticRegression(solver='lbfgs')
3 clf.fit(X, y)
4 print ("Nombre de bons prédictions avec Sickit Learn : ", np.sum(y == clf.predict(

```

Nombre de bons prédictions avec Sickit Learn : 100 sur 100

Pas mal du tout !

## Conclusion

**Attention dans les exemples précédents nous n'avons pas normalisé les données. Avant toute opération il est indispensable de le faire !**

Outre le fait qu'il puisse y avoir des données avec des valeurs très différentes, le fait de normaliser permet surtout de converger plus vite. Par exemple en utilisant la fonction suivante :

In [48]:

```

1 from sklearn import preprocessing
2
3 normalized_X = preprocessing.normalize(X)

```

Considérez la figure suivante. Précédemment nous avons vu qu'il fallait ajouter un 1 dans la première colonne de la matrice. Dans la figure ce  $x_0$  s'appelle le biais. Si nous appelons  $\sigma$ , la fonction de sortie, fonction d'activation, nous obtenons un perceptron qui est la base des algorithmes de réseaux de neurones ... et donc de l'apprentissage profond. La régression logistique ne vise pas le même objectif que les réseaux de neurones car il s'agit de déterminer dans le premier cas une probabilité sur la valeur de sortie et dans le second cas il constitue une étape vers d'autres couches éventuelles mais comme vous pouvez le constater la construction est très similaire, les  $\theta$  qui correspondent à nos paramètres sont généralement appelés *poids* dans les réseaux de neurones mais il s'agit de la même chose.

