

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

Descente de gradient

Mise en situation

Je travaille pour un site immobilier entre particuliers. Je souhaite fournir aux vendeurs une estimation raisonnable du prix de vente de manière automatisée.

Je dispose pour cela des données des dernières ventes, avec énormément de critères - probablement pas tous utiles...

Je n'arrive pas à mettre en place une régression linéaire : les milliers d'informations par bien vendu (sans parler de la montée de degré pour une régression polynomiale...) me donnent une matrice beaucoup trop grosse et trop longue à inverser :(

La descente de gradient

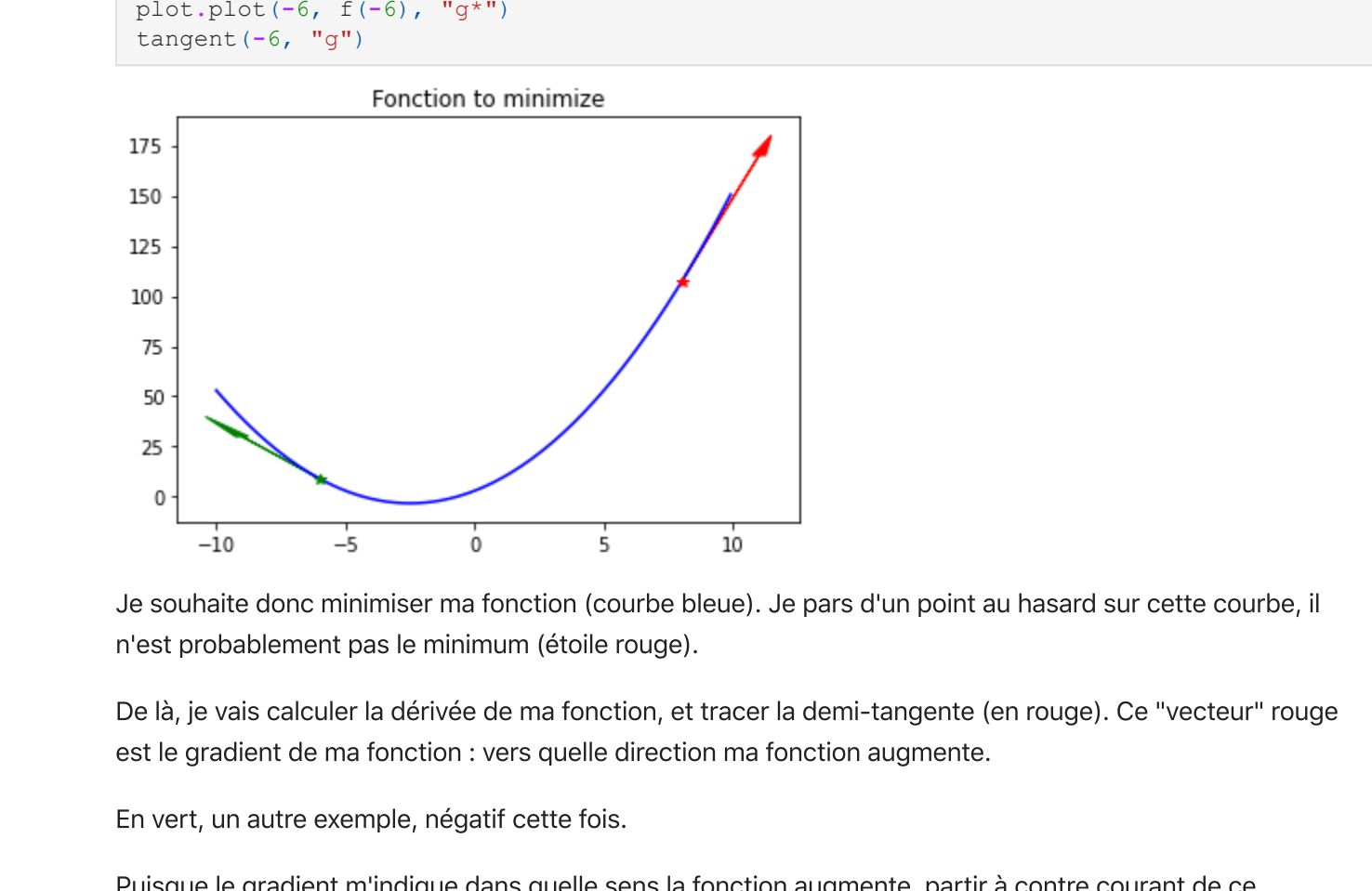
Le but de la méthode

Il est important de bien comprendre le problème ici. Toute la logique mathématique reste inchangée, il s'agit toujours de minimiser la fonction de coût par rapport aux paramètres du modèle.

Simplement, il faut trouver une autre méthode de minimisation. La résolution directe n'étant pas possible, il va falloir trouver un moyen de s'approcher du minimum, de manière rapide et précise de préférence. Et cette manière, c'est la descente de gradient.

Visualisation 2D

Pour mieux voir de quoi il retourne, un petit exemple visuel



Je souhaite donc minimiser ma fonction (courbe bleue). Je pars d'un point au hasard sur cette courbe, il n'est probablement pas le minimum (étoile rouge).

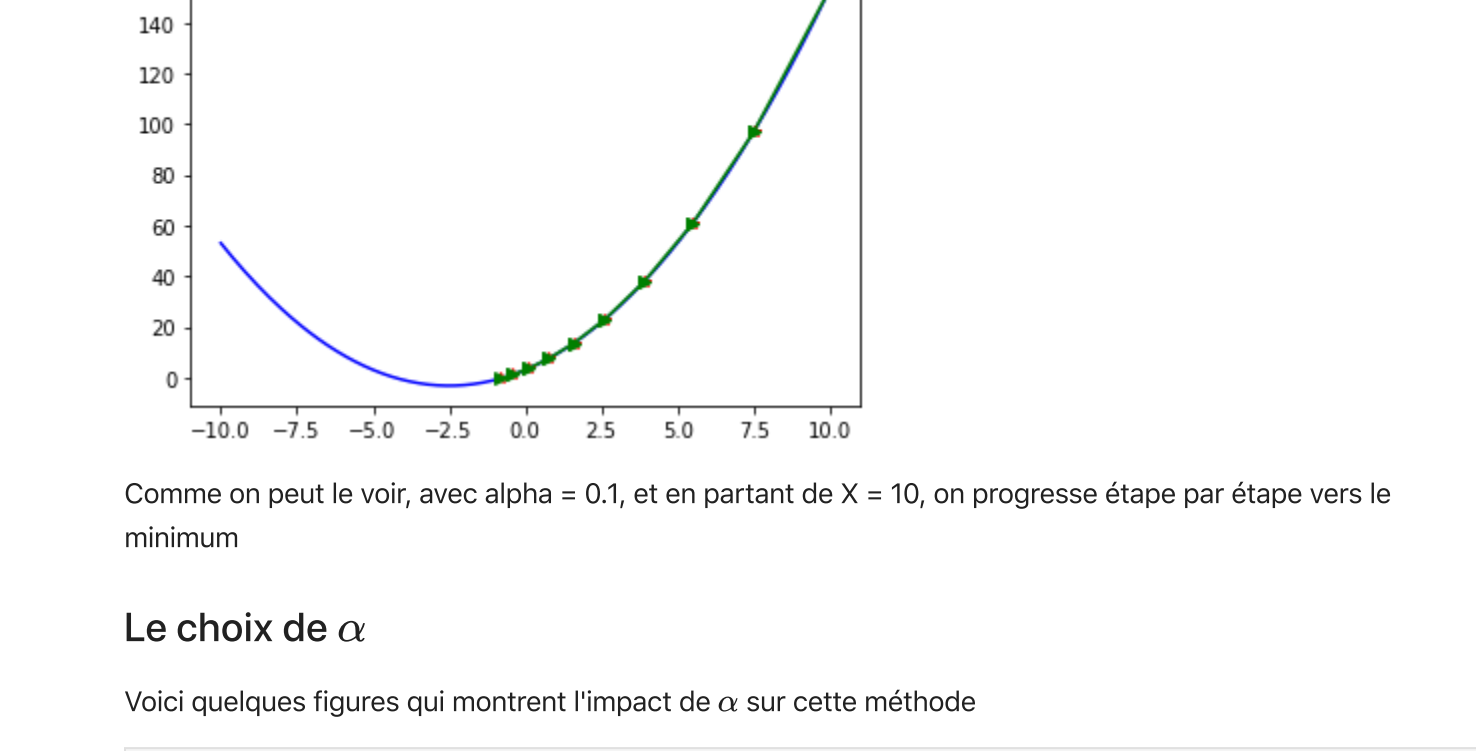
De là, je vais calculer la dérivée de ma fonction, et tracer la demi-tangente (en rouge). Ce "vecteur" rouge est le gradient de ma fonction : vers quelle direction ma fonction augmente.

En vert, un autre exemple, négatif cette fois.

Puisque le gradient m'indique dans quelle sens la fonction augmente, partir à contre courant de ce gradient devrait me rapprocher du minimum. Partant de là, on peut définir un processus itératif.

```
X = point au hasard
Répéter:
    dX = dérivée de f, en X
    X = X - alpha * dX
```

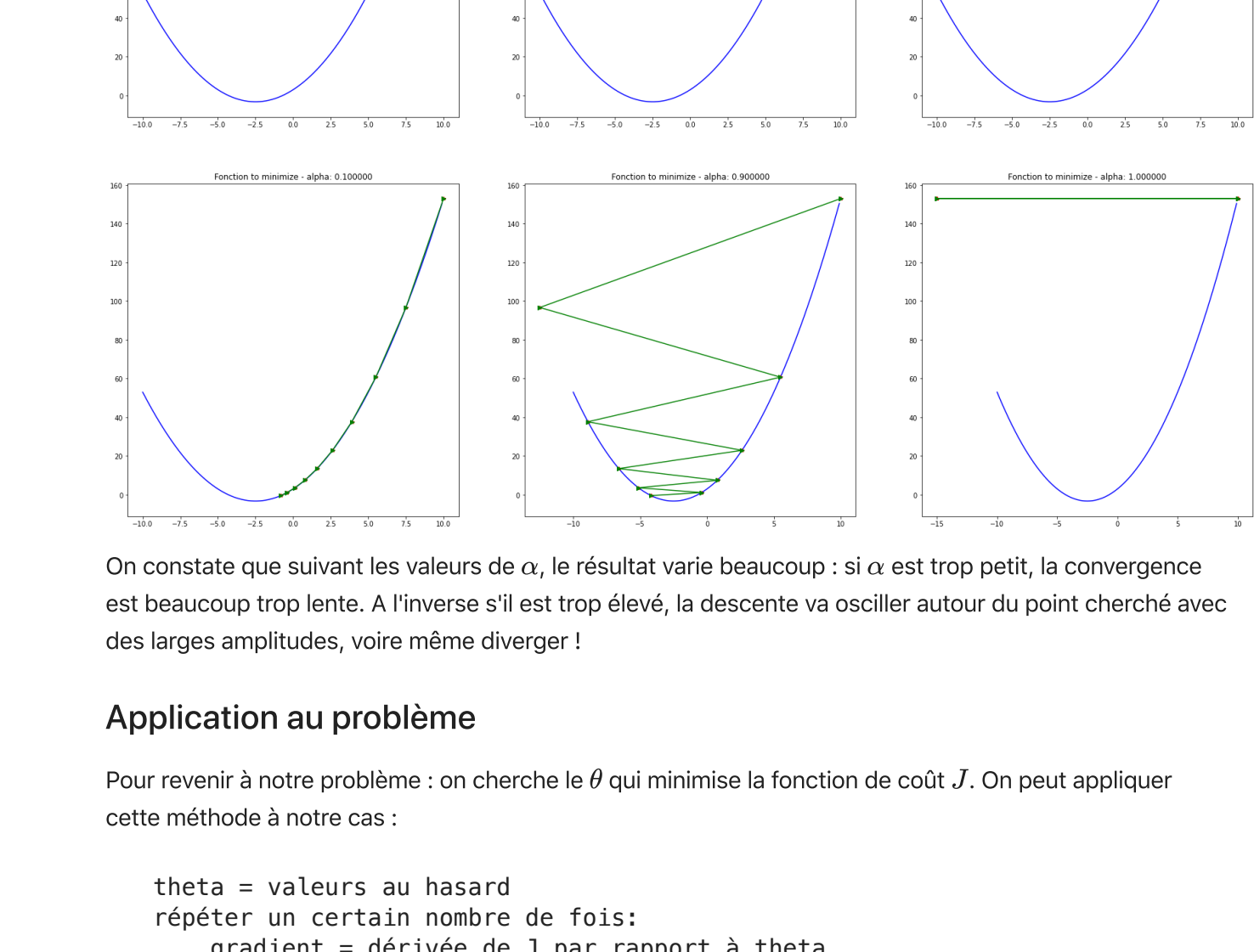
Notons l'introduction de alpha : le *learning rate*, ou taux d'apprentissage. On y reviendra rapidement, c'est un réel, on prendra 0.1 pour l'instant



Comme on peut le voir, avec alpha = 0.1, et en partant de X = 10, on progresse étape par étape vers le minimum

Le choix de α

Voici quelques figures qui montrent l'impact de α sur cette méthode



On constate que suivant les valeurs de α , le résultat varie beaucoup : si α est trop petit, la convergence est beaucoup trop lente. A l'inverse s'il est trop élevé, la descente va osciller autour du point cherché avec des larges amplitudes, voire même diverger !

Application au problème

Pour revenir à notre problème : on cherche le θ qui minimise la fonction de coût J . On peut appliquer cette méthode à notre cas :

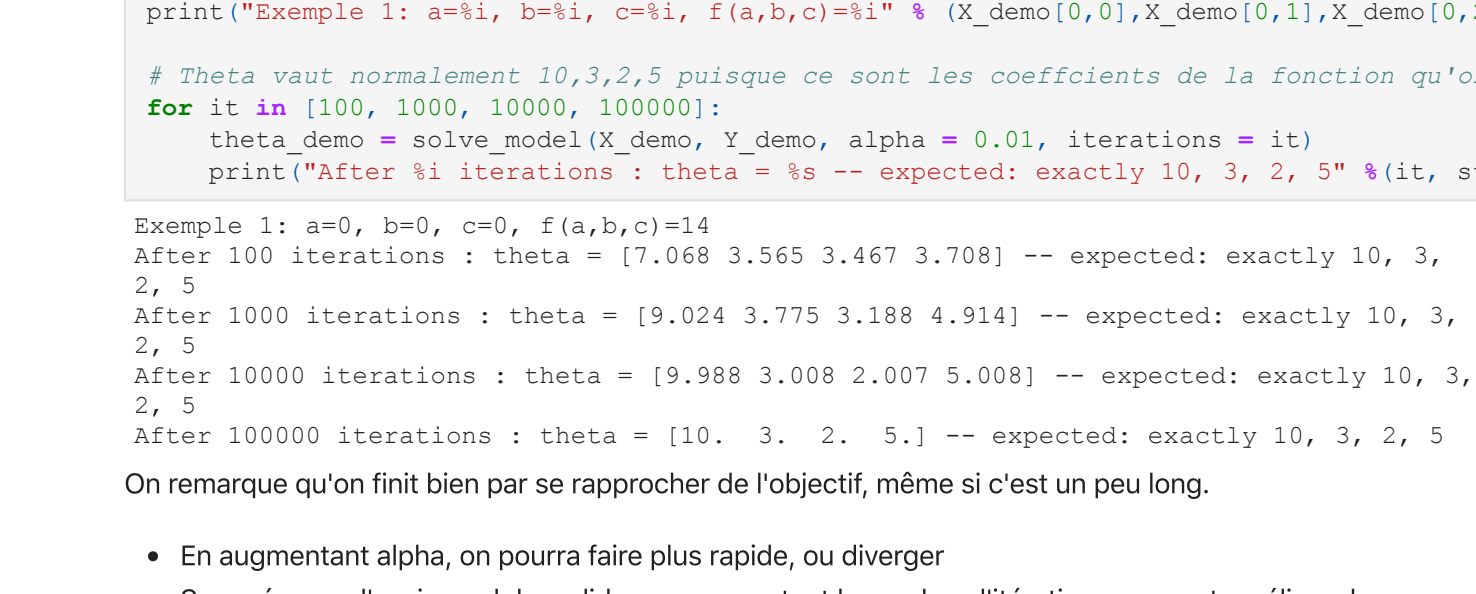
```
theta = valeurs au hasard
répéter un certain nombre de fois:
    gradient = dérivée de J par rapport à theta
    theta = theta - alpha * gradient
```

Et pour rappel, depuis le 1er jour, on a déjà la formule du gradient :

$$\frac{dJ}{d\theta} = \frac{1}{m} X^T \cdot (X \cdot \theta - Y)$$

Simulation rapide

Validons le concept avec une petite descente de gradient sur une fonction qu'on souhaite modéliser



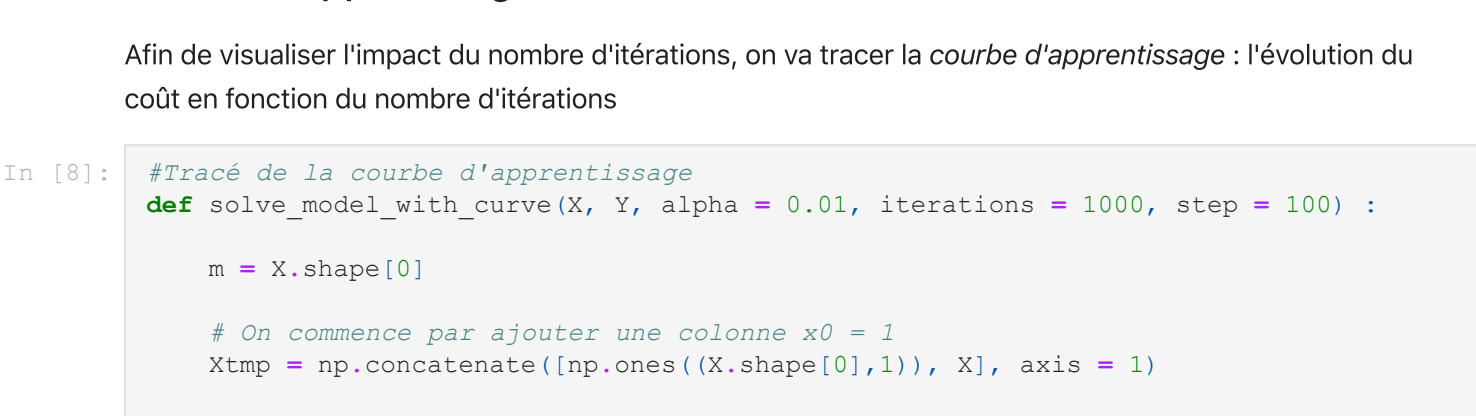
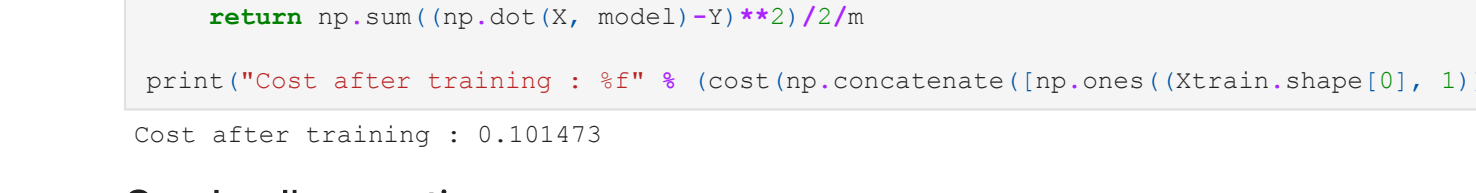
On remarque qu'on finit bien par se rapprocher de l'objectif, même si c'est un peu long.

- En augmentant alpha, on pourra faire plus rapide, ou diverger
- Sous réserve d'avoir un alpha valide, en augmentant le nombre d'itérations, on peut améliorer la précision

Retour à la mise en situation

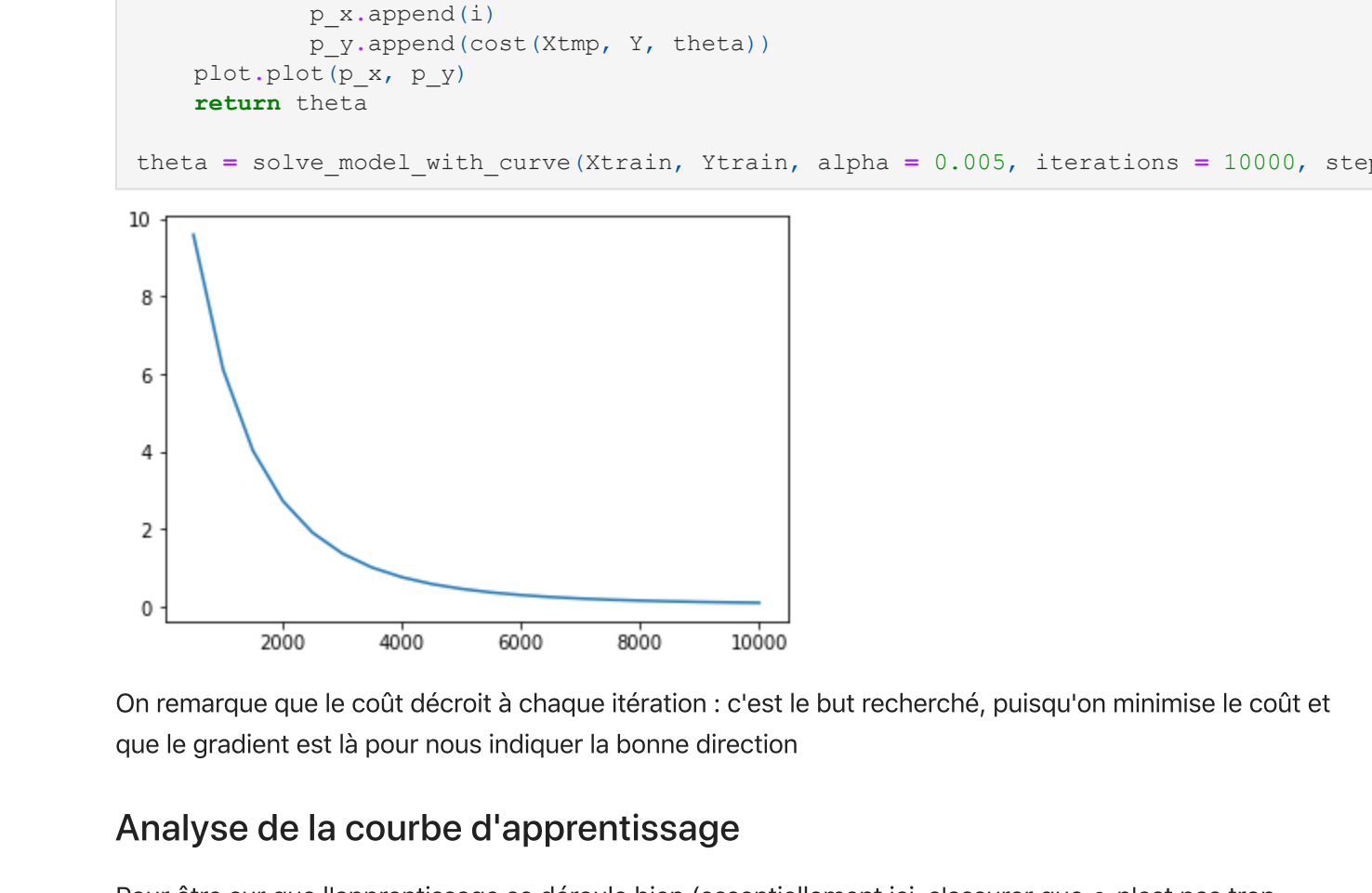
On va résoudre notre problème à l'aide de cette méthode

Chargement des données



Courbe d'apprentissage

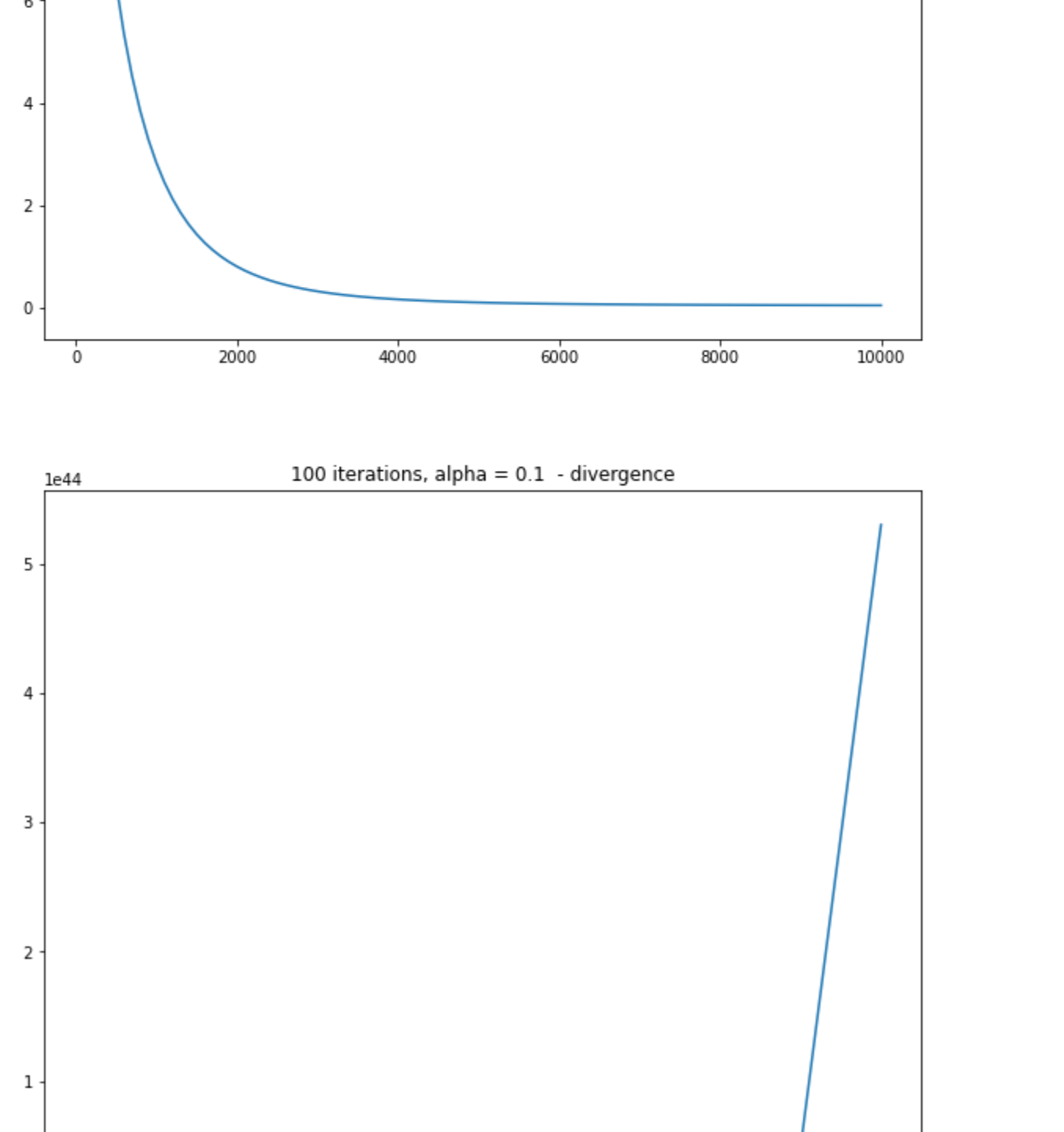
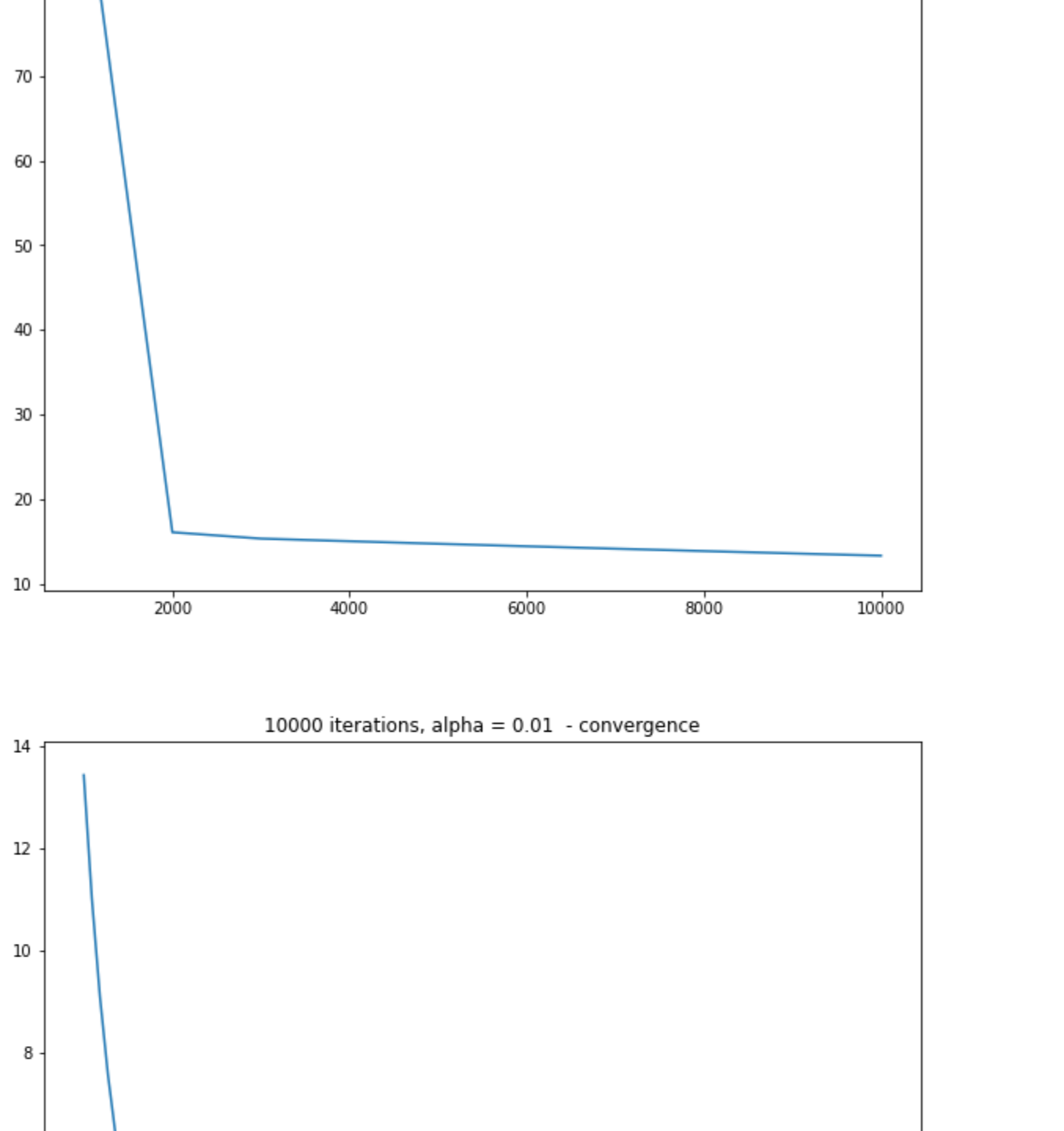
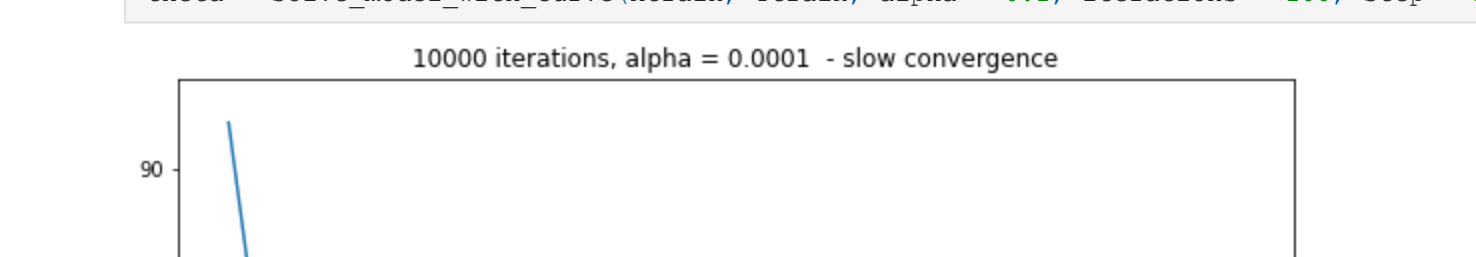
Afin de visualiser l'impact du nombre d'itérations, on va tracer la *courbe d'apprentissage* : l'évolution du coût en fonction du nombre d'itérations



On remarque que le coût décroît à chaque itération : c'est le but recherché, puisqu'on minimise le coût et que le gradient est là pour nous indiquer la bonne direction

Analyse de la courbe d'apprentissage

Pour être sur que l'apprentissage se déroule bien (essentiellement ici, s'assurer que α n'est pas trop grand, il est souvent intéressant d'analyser la courbe d'apprentissage justement.



- La 1ere courbe indique ce qui arrive avec un α trop petit. Ca converge, mais c'est très lent. Par contre, la courbe a quand même une bonne tête et on finira par converger
- La 2nde courbe présente l'apprentissage avec un "bon" α : la courbe décroît et se rapproche de zéro : le modèle est de plus en plus précis
- La dernière courbe est le cas d'un α beaucoup trop grand : le modèle diverge.

Paramètre vs hyperparamètre

Un paramètre du modèle est une des composantes qu'on fait varier pour minimiser la fonction de coût. Ex: θ

Un hyperparamètre du modèle est une valeur qu'on peut faire varier pour optimiser l'apprentissage. Ex: α

Un "mauvais" modèle est un modèle qui ne représentera jamais la fonction qu'on recherche. Exemple : essayer de modéliser une parabole avec une droite. Mais même pour un mauvais modèle, la fonction de coût admet un minimum et la descente de gradient doit y mener.

Pour résumer :

- si le modèle ne converge pas vers un minimum, c'est probablement la faute de α , qu'il faut faire varier
- sinon, le modèle converge forcément vers une valeur (qu'on définira plus tard comme le *biais* du modèle):
 - si cette valeur est proche de zéro, le modèle est probablement bon, en tous les cas il a trouvé ce qu'on lui demandait - mais était-ce la bonne chose à demander ? :-)
 - sinon, il faut changer de modèle, celui-ci ne trouvera probablement jamais la bonne représentation de la fonction recherchée