

Dans le cadre d'un sujet d'exercice à l'Ecole 42 Paris, je dois faire deux programmes :

Le premier est chargé de lire des données provenant d'un fichier csv, Ce fichier contenant des lignes avec deux informations (Kilométrage, et prix de voiture). Il doit ensuite utiliser une régression linéaire avec descente de gradient. (Ce projet assez simple nous donne déjà l'hypothèse à entraîner ainsi que le type de fonction que l'on doit trouver)

$$estimatePrice = \theta_0 + (\theta_1 * mileage)$$

il sauvegarde θ_0 et θ_1 pour être utilisé par le deuxième programme.

Celui-ci ne fait que récupérer ces infos pour que l'utilisateur rentre un kilométrage et que le programme lui donne l'estimation.

Mes programmes fonctionnent très bien et voici ma solution

Algo de Régression Linéaire

Hypothèse imposé :

les données sont organisées par : $\hat{y} = f(x) = \theta_0 + \theta_1 x$

pseudo code :

```
repeat until convergence* : {  
   $\nabla(J)$  compute  
     $\theta_0 := \theta_0 - \alpha \nabla(J)_0$   
     $\theta_1 := \theta_1 - \alpha \nabla(J)_1$   
}  
avec :  
 $\alpha$  learning rate (usually between 0 and 1)
```

*Pour l'instant on boucle un certain nombre de fois.

le gradient je le calcul comme ça :

$$\nabla(J) = \frac{1}{m} X'^T (X'\theta - y)$$

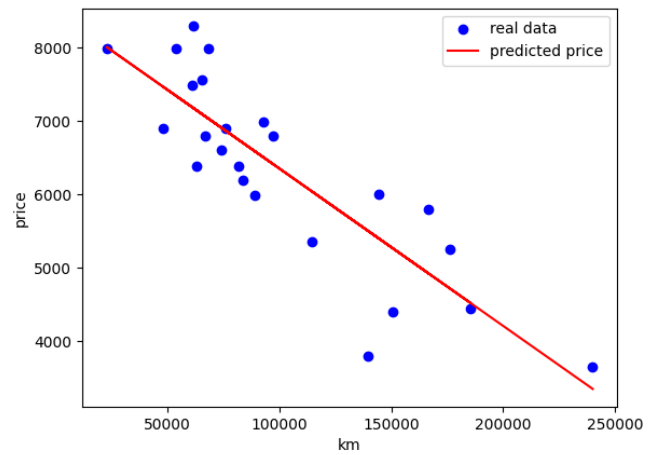
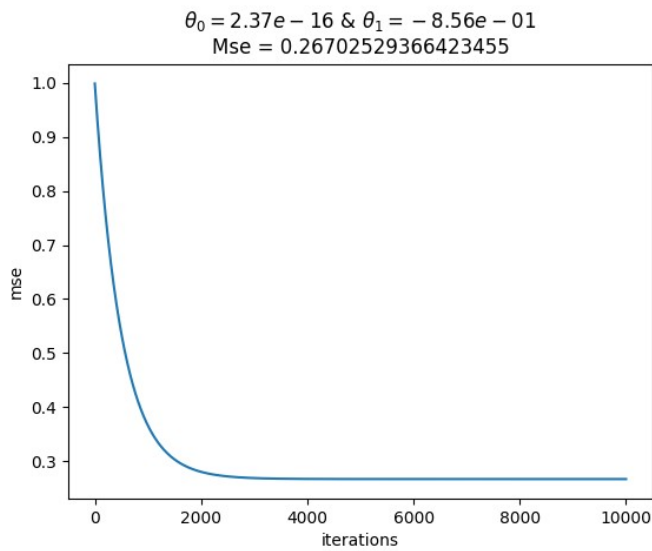
- $\nabla(J)$ vecteur de dimension 2 x 1.
- X' matrice (m x 2)
- X'^T matrice transposée de X' (2 x m)
- y vecteur de dimension m
- θ vecteur

en principe j'ai un objet **LinearRegression** qui a entre autre:

un attribut vecteur θ

une méthode **Predict(x)** qui calcul les prédictions avec les θ mémorisés.

Une méthode **fit()** qui : récupère la prédiction avec les θ , calcul le gradient ∇ (un tableau de deux float) et diminue θ_0, θ_1 avec le facteur α . il me renvoi un tableau des loss (MSE) pour afficher sa progression (On est sur une régression très simple et classique)



Ma question vient du fait que pour que je dérive correctement, j'ai dû Normaliser les données avec la formule

$$x'^i = \frac{x^i - \mu}{\sigma} \text{ avec } i \in [1, m]$$

x un vecteur de dimension m

x' la version normalisée du vecteur x

μ la moyenne des x

σ la variation standard

(Je normalise x et y)

donc mon programme 1 sauvegarde θ_0 et θ_1 'Normalisé' PLUS μ et σ des données x et y , pour que le deuxième programme fonctionne correctement.

Le programme 2 récupère le kilométrage(en km), le normalise et calcul un prix 'normalisé' par

$$\hat{y} = f(x) = \theta_0 + \theta_1 x$$

puis 'dé-normalise' pour donner à l'utilisateur un prix en euros...

et voici enfin ma question :

Peut on, dans le premier programme, calculer θ_0 et θ_1 qui seraient déjà 'dénormale' et ainsi ne pas enregistrer les moyennes et variations standards des vecteurs x , y ?