

TP

Apprentissage automatique

PSI : Lycée Rabelais

1 Présentation du problème

On s'intéresse à la prédiction du prix d'un bien immobilier en connaissant un certain nombre de critères. La base de données disponible a été établie dans les années 1990 dans le cadre de travaux de recherche. Elle s'intéresse aux biens immobiliers en Californie. Cette base est composée des caractéristiques suivantes en entrée :

- MedInc : le revenu médian des habitants du quartier ;
- HouseAge : l'âge médian des biens immobiliers du quartier ;
- AveRooms : le nombre de pièces dans le bien ;
- AveBedrms : le nombre de chambres à coucher ;
- Population : la population du quartier ;
- AveOccup : le nombre de personnes habitant le foyer ;
- Latitude : la latitude du quartier ;
- Longitude : la longitude du quartier.

La sortie est le prix médian du bien immobilier exprimé en centaine de milliers de dollars.

L'objectif est de définir un modèle permettant de prédire le prix du bien en fonction des caractéristiques préalablement mentionnées.

La trame du code est disponible sur votre *cahier-de-prépa* sous le nom `prix_immobilier__eleve.py`.

Question 1. Ouvrir, tester le code et commenter (ci-dessous) les principales instructions.

```
1 import numpy as np
2
3 from sklearn.datasets import fetch_california_housing
4 california_housing = fetch_california_housing()
5
6 from sklearn import datasets
7 prix = datasets.fetch_california_housing() ## importation de la bibliothèque
8
9 y = prix.target ## y contient le prix des biens immobiliers
10 x = prix.data ## X contient le tableau des données d'entrée
11
12 import seaborn as sns
13 import matplotlib.pyplot as plt
14
15 def trace_loc(x,y,T): ## cette fonction permet de tracer
16     ## le prix NORMALISÉ en fonction de la localité
17     ## NOTA : T permet d'afficher un commentaire
18
19     plt.figure()
20     ymap = (y - y.min()) / (y.max() - y.min())
21     sns.scatterplot(x[:, -2], x[:, -1], ymap, cmap="viridis")
22     plt.xlabel('Longitude (en degrés)')
23     plt.ylabel('Latitude (en degrés)')
24     plt.title("Prix normalisé de l'immobilier (1 = prix max ; 0 = prix min)"+'\n'+T)
25     plt.grid(True)
```

```

26 def trace_simple(x,i,y,T): ## cette fonction permet de tracer
27     ## le prix en fonction de la i-ème caractéristique
28     ## NOTA : T permet d'afficher un commentaire
29
30     plt.figure()
31     plt.plot(x[:,i],y,'.')
32     features = ["du revenu médian des habitants du quartier","de l'age médian des biens
33     immobiliers du quartier","du nombre de pièces dans le bien","du nombre de chambres à
34     coucher","de la population du quartier","du nombre de personnes habitant le foyer","de la
35     latitude du quartier","de la longitude du quartier"]
36     plt.xlabel('En fonction '+features[i])
37     plt.ylabel('Prix en centaine de milliers de dollars')
38     plt.title("Prix de l'immobilier en fonction "+features[i]+'\\n'+T)
39     plt.grid(True)
40
41 #####
42 ##### création du modèle de régression et exportation des résultats
43 #####
44 def apprentissage():
45     ## pour importer MLPRegressor
46     from sklearn.neural_network import MLPRegressor
47     mlp = MLPRegressor(activation='identity',hidden_layer_sizes=(4,2)) ## Création du modèle
48
49     from sklearn.model_selection import train_test_split
50     ## Préparation des données
51     x_train, x_test, y_train, y_test = train_test_split(x, y,shuffle=True, test_size=0.5)
52     ## nota : test_size=0.5 permet d'avoir autant de données de test que de données d'
53     apprentissage
54
55     mlp.fit(x_train,y_train)
56     y_pred = mlp.predict(x_test)
57
58     return(x_train, x_test, y_train, y_test,y_pred,mlp) ## Score obtenu
59
60 #####
61 ##### récupération des résultats
62 #####
63 x_train, x_test, y_train, y_test,y_pred,mlp = apprentissage()
64 score = mlp.score(x_test, y_test)
65 trace_loc(x_train,y_train,"!!!! Tracé des données d'apprentissage !!!!")
66 trace_loc(x_test,y_pred,"!!!! Tracé des données prédites !!!!")
67 trace_simple(x_train,0,y_train,"!!!! Tracé des données d'apprentissage !!!!")
68 print("prix minimal prédit",min(y_pred))
69 print("prix maximal prédit",max(y_pred))
70 print("prix poyen prédit",np.mean(y_pred))

```

Question 2. Combien de paramètres (poids et biais) sont présents dans le modèle ? Le nombre de données est-il pertinent ?

Question 3. Utiliser la fonction trace_simple (sans relancer entièrement la phase d'apprentissage) pour visualiser l'influence de chacune des caractéristiques d'entrée sur le prix de l'immobilier.

Question 4. Donner une unité à chacune des caractéristiques d'entrée.

Un problème récurrent concerne la gestion des entrées lorsque celles-ci ne représentent pas la même grandeur.

On peut constater qu'une entrée joue un rôle prépondérant, si celle-ci est, en valeur numérique, plus grande que les

autres. Pour éviter ce problème, une étape fréquemment utilisée est donc la normalisation des données. Il s'agit, par exemple, de calculer pour une entrée x :

$$x_{\text{normalisée}} = \frac{x}{x_{\text{max}}}$$

Ce problème étant fréquent, des modules adaptés permettent la normalisation de tout le tableau des entrées.

Question 5. Rajouter les lignes suivantes, permettant la normalisation, avant la phase d'apprentissage et observer l'effet sur le score de l'algorithme.

```
1 from sklearn.preprocessing import StandardScaler
2 scaler = StandardScaler()
3 x_train = scaler.fit_transform(x_train)
4 x_test = scaler.fit_transform(x_test)
```

Question 6. Commenter les nuages de points obtenus et notamment :

- d'un point de vue de la géographie (on pourra chercher une carte de la Californie sur internet si besoin).
- en comparant les données d'apprentissage et celles prédites.

Question 7. Écrire une fonction permettant de calculer le score moyen sur n_{essai} essais.

Question 8. Commenter les résultats obtenues d'un point de vue du score et en visualisant les résultats. Les prédictions semblent-elles pertinentes ?

La documentation indique :

 score : Return the coefficient of determination R^2 of the prediction.

The coefficient of determination R^2 is defined as $1 - \frac{u}{v}$, where u is the residual sum of squares $((y_{\text{true}} - y_{\text{pred}})**2).sum()$ and v is the total sum of squares $((y_{\text{true}} - y_{\text{true}.mean()})**2).sum()$. The best possible score is 1.0 and it can be negative (because the model can be arbitrarily worse). A constant model that always predicts the expected value of y , disregarding the input features, would get a score of 0.0.

On peut également évaluer l'intérêt du réglage du taux d'apprentissage (*learning rate*) en rajoutant dans le modèle :

```
mlp = MLPRegressor(activation='identity',hidden_layer_sizes=(10,10),learning_rate_init=0.001)
                (ici on impose un taux d'apprentissage de 0.001)
```

On peut également afficher le temps d'apprentissage des paramètres en modifiant dans le code précédent :

```
1 from time import time
2 start = time()
3 apprentissage()
4 t_app = time()- start
```

Question 9. Faire quelques simulations (environ 5 pour chaque taux d'apprentissage) et remplir le tableau ci-dessous puis conclure.

Taux d'apprentissage	Temps d'apprentissage moyen (en s)	Score moyen
0.001		
0.01		
0.1		
0.5		

Question 10. Le score est-il toujours un indicateur objectif ?

Question 11. Modifier le code pour que celui-ci ne prenne en compte que la latitude et la longitude du bien immobilier (on rappelle que, pour un tableau numpy de dimension $n \times m$ noté TAB, TAB[:, i:j] renvoie toutes les lignes et les colonnes allant de la i -ième (inclusive) à la j -ième exclusive).

Question 12. Observer le résultat puis commenter les lignes rajoutées pour ne plus les prendre en compte dans la suite.

2 Amélioration du modèle

Dans le modèle préalablement défini, on a choisi une fonction identité pour la fonction d'activation.

Question 13. Modifier cette fonction d'activation en imposant une fonction sigmoïde (logistic). Observer le score et conclure.

Question 14. Modifier la structure du réseau en imposant le nombre de neurones sur chaque couche cachée en écrivant, par exemple, hidden_layer_sizes=(40,20,10) puis ensuite hidden_layer_sizes=(2). Cela signifierait, dans le premier exemple, qu'il y a trois couches cachées de 40 puis 20 puis 10 neurones. Déterminer le nombre de paramètres associés à chaque configuration. Observer l'effet sur le score obtenu.

Question 15. Quelle est la condition principale pour avoir un modèle très performant ?

Question 16. Définir le réseau le plus performant. **Attention, votre réseau ne doit pas surapprendre !**

3 Utilisation d'un algorithme des k -plus proches voisins

Question 17. Rappeler le fonctionnement d'un tel algorithme.

Question 18. Modifier le modèle en écrivant :

```
1 | from sklearn.neighbors import KNeighborsRegressor  
2 | mlp = KNeighborsRegressor(n_neighbors=3) ## Création du modèle
```

Question 19. Observer le score. Choisir le nombre de voisins pour que votre modèle soit le plus performant possible.