

TP introductif en Big Data

Saïda Guezoui, groupe DS

Octobre 2021

1 Introduction

Le Big Data prend aujourd'hui une place importante dans le marché des entreprises et il est de plus en plus utilisé.

Le bootstrap est un algorithme permettant l'obtention de plusieurs échantillons à partir de notre échantillon origine. Cet algorithme consiste à faire un rééchantillonnage avec remise d'un jeu de donnée pour avoir des échantillons statistiques qui ont la même taille que notre échantillon initial.

Il existe deux types de bootstrap, paramétrique et non paramétrique.

Dans ce TP, nous étudierons ces différents bootstrap à l'aide des bibliothèques et modules en Python suivant

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 from sklearn.linear_model import LogisticRegression
5 from sklearn.linear_model import LinearRegression
6 import seaborn as sns
```

Listing 1: Modules Python3

2 Etude du jeu de donnée "*student_score*"

2.1 Importation du jeu de données et calcul de la matrice de corrélation

```
1 # Importation des données du fichier txt
2 score = np.genfromtxt(fname = "student_score.txt")
3
4 # Creation d'un DataFrame
5 score_student = pd.DataFrame(score, columns = ['mech' , 'vecs', 'alg', 'analy', 'stat'])
6
7 # On supprime la première ligne qui ne contient que des NaN
8 score_student = score_student.drop(data.index[0])
9 print("Données : \n", score_student)
10
11 # Calcul des corrélations entre les variables
12 correlation = score_student.corr()
13 print("\n Matrice de corrélation :")
14 sns.heatmap(data = correlation, annot=True)
```

2.2 Estimation du paramètre θ

```
1 # Calcul des valeurs et vecteurs propres de la matrice des correlations
2 v, vp = np.linalg.eig(correlation)
3
4 # Calcul de Theta
5 theta = max(v)/ sum(v)
6 print("\n Theta = ", theta)
```

La valeur de θ est d'environ 0.7.

2.3 Définition de la fonction bootstrap

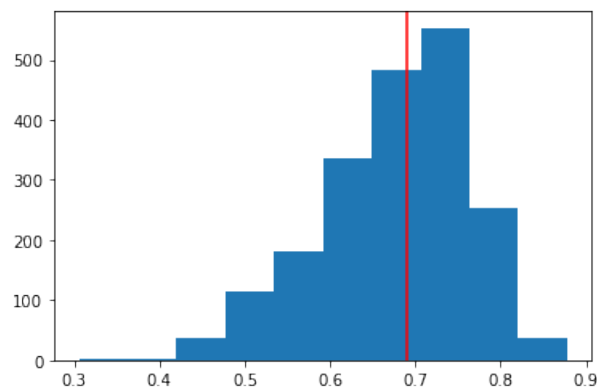
Nous définissons la fonction bootstrap comme une fonction qui prend en paramètre un data donné et nombre B fixé à 2000 dans cette partie. On cherche à obtenir les valeurs de notre vecteur $\hat{\theta}$, nous faisons donc un rééchantillage avec remise sur notre jeu de données, puis sur la matrice de correlation. Ensuite, nous appliquons la formule de θ suivante pour avoir notre liste des $\hat{\theta}$:

$$\theta = \frac{\lambda_{max}}{\sum(\lambda)}$$

Avec λ_{max} est la valeur propre maximale. En appliquant cette fonction au jeu de données "score_student", nous obtenons alors un vecteur de $\hat{\theta}$ de taille 2000.

```
1 def bootstrap(data,B):
2     theta_list = []
3
4     for _ in range(B):
5         # Echantillon tire et replace = True signifie un tirage avec remise
6         sample = data.sample(n = len(data.index), replace = True)
7
8         correlation = sample.corr()
9
10        eigenvalue, eigenvector = np.linalg.eig(correlation)
11
12        theta = max(eigenvalue)/sum(eigenvalue)
13        theta_list.append(theta)
14
15    return theta_list
16
17 theta_list = bootstrap(score_student,B = 2000)
```

2.4 Représentation sous forme d'histogramme



2.5 L'erreur quadratique moyenne

L'erreur quadratique moyenne est d'environ 0.086.

3 Etude du jeu de données "diabete"

Dans cette partie, on s'intéresse au jeu de données "diabete" contenant des informations sur diabète des indiennes Pima. Nous cherchons à prédire la survenue du diabète à l'aide de la régression logistique en utilisant les autres covariables.

3.1 Importation du jeu de données

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	Diabete	Age	Outcome
0	6	148	72	35	0	33.6	0.627	50	1
1	1	85	66	29	0	26.6	0.351	31	0
2	8	183	64	0	0	23.3	0.672	32	1
3	1	89	66	23	94	28.1	0.167	21	0
4	0	137	40	35	168	43.1	2.288	33	1

3.2 Création du Design, de la variable à prédire "Outcome"

3.3 Application du modèle de regression logistique

```
1 # Creation du modele
2 model_pima = LogisticRegression(max_iter = 300)
3
4 # Apprentissage
5 model_pima.fit(Design_pima, Y_pima)
6
7 # Test du modele (score)
8 print("Score : ", model_pima.score(Design_pima, Y_pima))
9
10 # Parametres du modele
11 intercept_pima = model_pima.intercept_
12 coefs_pima = model_pima.coef_
```

Définition de la fonction logit et du p .

```
1 # Logit
2 def g(x):
3     return np.exp(x) / (1 + np.exp(x))
4
5 # Calcul de p_hat
6 p_hat = g(intercept_pima + coefs_pima.dot(Design_pima.T))[0]
```

3.4 Application du bootstrap paramétrique

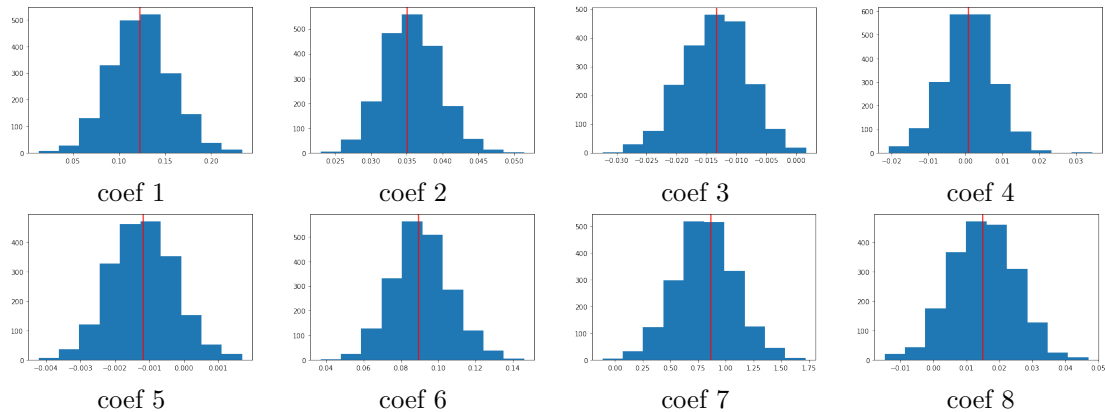
```
1 def bootstrap_para(B) :
2     Coefs = []
3     Intercepts = []
4
5     for _ in range(B) :
6         Y_pima_hat = np.random.binomial(n = 1, p = p_hat, size = len(p_hat))
7         model_pima.fit(Design_pima, Y_pima_hat)
8         Coefs.append(model_pima.coef_[0])
9         Intercepts.append(float(model_pima.intercept_))
10
11     return Coefs, Intercepts
```

```

12
13 vect_coef, liste_intercepts = bootstrap_para(2000)

```

Affichage des distributions des coefficients sous forme d'histogrammes



4 Etude du jeu de données "housing"

Nous intéressons cette fois-ci aux prix des logements dans différents quartiers de Boston. La variable en question est "MEDV". Nous cherchons à prédire ce prix en fonction des autres covariables lstat, ptratio et rm en utilisant un modèle de régression linéaire et un algorithme de bootstrap non paramétrique.

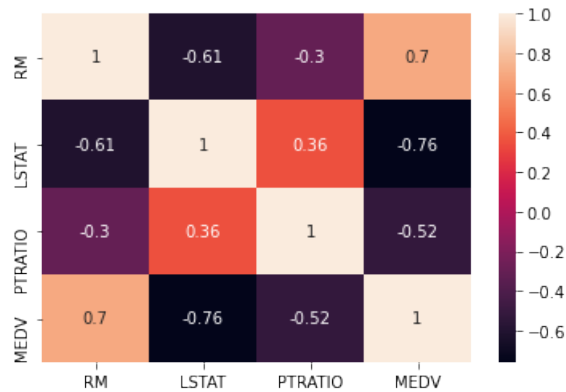
4.1 Importation du jeu de données housing et calcul de la matrice de corrélation

```

1 housing = pd.read_csv("housing.csv")
2 housing = pd.DataFrame(housing, columns = ['RM', 'LSTAT', 'PTRATIO', 'MEDV'])
3 print(housing)
4
5 # Suppression de la premiere colonne d'indice
6 housing = housing.drop(housing.index[0])
7
8 # Etude de la correlation
9 correlation = housing.corr().round(2)
10 sns.heatmap(data = correlation, annot=True)

```

	RM	LSTAT	PTRATIO	MEDV
0	6.575	4.98	15.3	504000.0
1	6.421	9.14	17.8	453600.0
2	7.185	4.03	17.8	728700.0
3	6.998	2.94	18.7	701400.0
4	7.147	5.33	18.7	760200.0



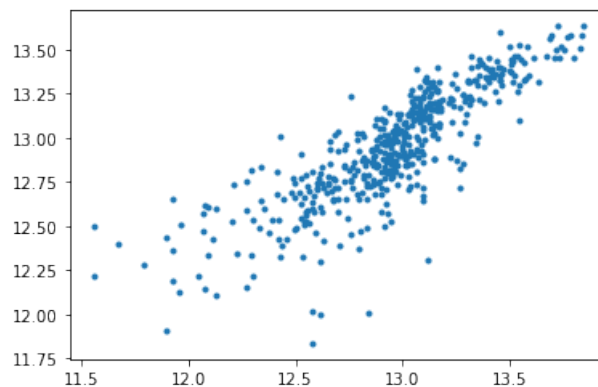
Nous constatons que la variable MEDV est corélee positivement avec la variable RM et négativement avec les variables LSTAT et PTRATIO. Nous pouvons donc appliquer un modèle de régression linéaire.

4.2 Création du Design et de la variable à prédire $\log(\text{MEDV})$

```
1 # Variable a predire Y = log(MEDV)
2 Y_housing = np.log(housing['MEDV'])
3 # Design
4 Design_housing = housing.drop('MEDV',1)
```

4.3 Application de la régression linéaire

```
1 REG_housing = LinearRegression()
2 # Entraînement du modele
3 Results = REG_housing.fit(Design_housing, Y_housing)
4
5 # Coefs et intercept
6 Coefs_housing = REG_housing.coef_
7 Intercept_housing = REG_housing.intercept_
8
9 # Affichage
10 print("Les coefficients de la regression lineaire : \n", Coefs_housing)
11 print("\nL'intercept de la regression lineaire : \n", Intercept_housing)
12
13 # Nuage de points de la regression
14 plt.plot(Y_housing, REG_housing.predict(Design_housing),'.')
15 plt.show()
```



4.4 Calcul du σ^2 et ϵ

```

1 RMSE = np.sqrt(((Y_housing - REG_housing.predict(Design_housing))*2).sum()/len(
    Y_housing))
2 print(RMSE)
3
4 Epsilon = Y_housing - Intercept_housing - Coefs_housing.dot(Design_housing.T)
5 print(Epsilon)

```

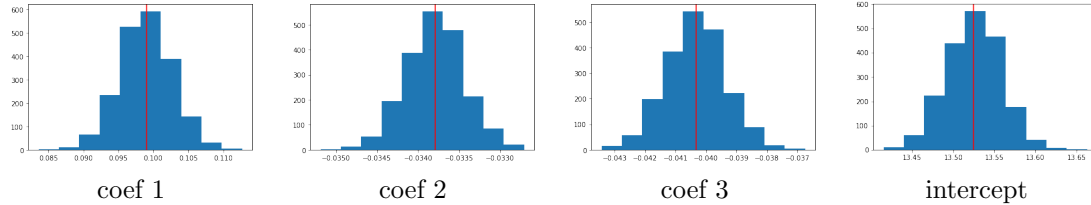
4.5 Application du bootstrap non paramétrique

```

1 def bootstrap_non_para(B) :
2     Coefs = []
3     Intercepts = []
4
5     for _ in range(B) :
6         Epsilon_star = np.random.normal(0, RMSE**2, len(Epsilon))
7         Y_housing_hat = Intercept_housing + Coefs_housing.dot(Design_housing.T) +
            Epsilon_star
8         REG_housing.fit(Design_housing, Y_housing_hat)
9         Coefs.append(REG_housing.coef_)
10        Intercepts.append(float(REG_housing.intercept_))
11
12    return Coefs, Intercepts
13
14 Coefs_h, Intercepts_h = bootstrap_non_para(2000)

```

Représentation des distributions des intercepts et coefs sous forme d'histogrammes



```

1 def bootstrap(data,B):
2     coefs = []
3     intercepts = []
4
5     REG_housing = LinearRegression()
6
7     for _ in range(B):
8         sample = data.sample(n = len(data.index), replace = True)
9
10        X = sample.drop('MEDV',1)
11        y = np.log(sample['MEDV'])
12
13        REG_housing.fit(X, y)
14
15        coefs.append(REG_housing.coef_)
16        intercepts.append(REG_housing.intercept_)
17
18    return coefs, intercepts
19
20 coefs_housing, intercepts_housing = bootstrap(housing,B = 2000)

```

Représentation des distributions des intercepts et coefs sous forme d'histogrammes

