



Session d'octobre 2024

Table des matières

1	Introduction	2
2	Introduction aux bibliothèques utilisees	2
3	Partie 1 : Visualisation des données	3
3.1	Exploration des jeux de données	3
3.2	Resumé statistique	5
3.3	Nuage des points du premier data set « Experience / Salaire »	6
3.4	Nuages des points du deuxième data set selon les proprietes « Features »	7
4	Partie 2 : Regression lineaire simple (Experience vs Salaire)	8
4.1	Entraînement de modèle	8
4.2	Prediction sur le dataset de test	9
4.3	Visualisation des resultats de la regression	9
4.4	Evaluation du model	10
4.4.1	Interpretation	10
5	Partie 3 : Regression multiple (Assurance)	11
5.1	EDA (Exploratory Data Analysis)	11
5.2	Selectionner 3 proprietes selon leurs degre d'importance	14
5.3	Technique de standardisation	15
5.4	Entraîner le modèle	16
5.5	Prediction	17
5.6	Visualisation	17
5.7	Evaluation de modèle en utilisant ces trois methodes	18
5.7.1	Interpretation	19
6	Partie 4 : Regression polynomiale (China GDP)	20
6.1	Comparaison des modèles	20
6.2	Prediction des données d'un data set de test pour les deux modèles	21
6.3	Visualisation des resultats des deux modèles	21
6.4	Evaluation des deux modèles en utilisant ces trois methodes	23
6.4.1	Interpretation	24
7	Conclusion	25
8	References	25

1 Introduction

L'objectif de cet atelier est de comprendre et de mettre en pratique les concepts fondamentaux de la regression lineaire simple et multiple, ainsi que la regression polynomiale, en s'appuyant sur des données reelles. Ce rapport documente les etapes, les analyses, et les conclusions tirees de ces experiences. La regression est un outil cle dans l'apprentissage automatique, permettant de modeliser des relations entre des variables et de faire des predictions precises. Nous avons choisi des datasets varies pour mieux comprendre l'impact des modèles et techniques sur des données reelles, en explorant differentes approches analytiques et de visualisation.

2 Introduction aux bibliothèques utilisees

Pour mener à bien cette etude, plusieurs bibliothèques Python essentielles ont ete importees.

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from pandas.plotting import scatter_matrix
5 from sklearn.model_selection import train_test_split
6 from sklearn.linear_model import LinearRegression
7 from sklearn.preprocessing import PolynomialFeatures, StandardScaler
8 from sklearn.metrics import mean_squared_error, mean_absolute_error,
   r2_score
9 import seaborn as sns
```

Voici une brève description de chacune d'elles :

- **pandas** : Utilisee pour la manipulation et l'analyse des données, notamment pour lire et traiter les ensembles de données sous forme de DataFrame.
- **numpy** : Fournit des fonctions mathematiques avancees pour les calculs numeriques et est souvent utilise pour travailler avec des tableaux multidimensionnels.
- **matplotlib.pyplot** : Une bibliothèque de visualisation pour generer des graphiques simples ou complexes.
- **pandas.plotting.scatter_matrix** : Un outil de pandas permettant de generer des matrices de graphiques de dispersion pour analyser les relations entre plusieurs variables.
- **sklearn.model_selection.train_test_split** : Permet de diviser les ensembles de données en ensembles d'entraînement et de test pour la modelisation.
- **sklearn.linear_model.LinearRegression** : Implemente des modèles de regression lineaire pour la prediction des valeurs continues.

- **sklearn.preprocessing.PolynomialFeatures** et **StandardScaler** : Utilisées respectivement pour transformer les caractéristiques en polynômes et pour standardiser les données en les centrant et en les réduisant.
- **sklearn.metrics.mean_squared_error**, **mean_absolute_error** et **r2_score** : Fournissent des outils pour évaluer la performance des modèles, notamment à l'aide des erreurs quadratiques moyennes, des erreurs absolues moyennes et du coefficient de détermination R^2 .
- **seaborn** : Une bibliothèque de visualisation statistique basée sur **matplotlib**, facilitant la création de graphiques plus esthétiques et informatifs.

Ces outils combinés permettent une exploration approfondie des données, une modélisation performante et une évaluation robuste des résultats obtenus.

3 Partie 1 : Visualisation des données

3.1 Exploration des jeux de données

Les jeux de données utilisés sont les suivants :

- **Experience et Salaire** : Ce dataset relie les années d'expérience professionnelle à un salaire annuel.
- **Assurance** : Ce dataset inclut des données démographiques et des informations sur les coûts d'assurance médicale.
- **China GDP** : Ce dataset montre l'évolution du PIB chinois sur plusieurs décennies.

Le code ci-dessous permet de charger et d'explorer les données :

```

1 # Chargement des datasets
2 experience_salary = pd.read_csv("./Salary_Data.csv")
3 insurance = pd.read_csv("./insurance.csv")
4 china_gdp = pd.read_csv("./china_gdp.csv")
5
6
7 # Partie 1 : Exploration et visualisation
8 # Question 1 : Explorer les données
9 # Exploration des données des deux datasets
10 print("----- Exploration du dataset Experience-Salaire :
    -----")
11 print(experience_salary.info())
12 print(experience_salary.head())

```

```

13
14 print("\n----- Exploration du dataset Assurance : -----")
15 print(insurance.info())
16 print(insurance.head())

```

Les resultats de cette etape sont illustres dans la Figure 1.

```

----- Exploration du dataset Expérience-Salaire : -----
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 30 entries, 0 to 29
Data columns (total 2 columns):
#   Column          Non-Null Count  Dtype
---  ---
0   YearsExperience  30 non-null    float64
1   Salary           30 non-null    float64
dtypes: float64(2)
memory usage: 612.0 bytes
None
   YearsExperience  Salary
0              1.1  39343.0
1              1.3  46205.0
2              1.5  37731.0
3              2.0  43525.0
4              2.2  39891.0

----- Exploration du dataset Assurance : -----
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1338 entries, 0 to 1337
Data columns (total 7 columns):
#   Column      Non-Null Count  Dtype
---  ---
0   age         1338 non-null  int64
1   sex         1338 non-null  object
2   bmi         1338 non-null  float64
3   children    1338 non-null  int64
4   smoker      1338 non-null  object
5   region      1338 non-null  object
6   charges     1338 non-null  float64
dtypes: float64(2), int64(2), object(3)
memory usage: 73.3+ KB
None
   age  sex  bmi  children  smoker  region  charges
0   19  female  27.900      0    yes  southwest  16884.92400
1   18   male  33.770      1    no   southeast  1725.55230
2   28   male  33.000      3    no   southeast  4449.46200
3   33   male  22.705      0    no  northwest  21984.47061
4   32   male  28.880      0    no  northwest  3866.85520

```

FIGURE 1 – Exploration des jeux de données.

3.2 Résumé statistique

Pour mieux comprendre la structure des données, voici le Résumé statistique :

```
1 # Question 2 : Resum statistique avec interpretation
2 # Resum statistique pour le dataset Experience-Salaire
3 print("----- Resum statistique du dataset Experience-Salaire :
  -----")
4 print(experience_salary.describe())
5
6 # Resum statistique pour le dataset Assurance
7 print("\n----- Resum statistique du dataset Assurance :
  -----")
8 print(insurance.describe())
```

Les resultats de cette etape sont illustres dans la Figure 2.

```
----- Résumé statistique du dataset Expérience-Salaire : -----
      YearsExperience      Salary
count      30.000000      30.000000
mean         5.313333    76003.000000
std          2.837888    27414.429785
min          1.100000    37731.000000
25%          3.200000    56720.750000
50%          4.700000    65237.000000
75%          7.700000   100544.750000
max         10.500000   122391.000000

----- Résumé statistique du dataset Assurance : -----
      age      bmi      children      charges
count  1338.000000  1338.000000  1338.000000  1338.000000
mean    39.207025    30.663397    1.094918   13270.422265
std     14.049960     6.098187    1.205493   12110.011237
min     18.000000    15.960000    0.000000    1121.873900
25%     27.000000    26.296250    0.000000    4740.287150
50%     39.000000    30.400000    1.000000    9382.033000
75%     51.000000    34.693750    2.000000   16639.912515
max     64.000000    53.130000    5.000000   63770.428010
```

FIGURE 2 – Résumé statistique des jeux de données.

3.3 Nuage des points du premier data set « Experience / Salaire »

Voici les visualisations demandées :

```
1 # Question 3 : Nuage des points pour le dataset Experience-Salaire
2 experience_salary.plot(kind='scatter', x='YearsExperience', y='Salary',
3     title='Experience vs Salaire')
4
5 plt.xlabel('Annees d\'experience')
6 plt.ylabel('Salaire')
7 plt.title('Nuage de points : Experience vs Salaire')
8 plt.show()
9 plt.show()
```

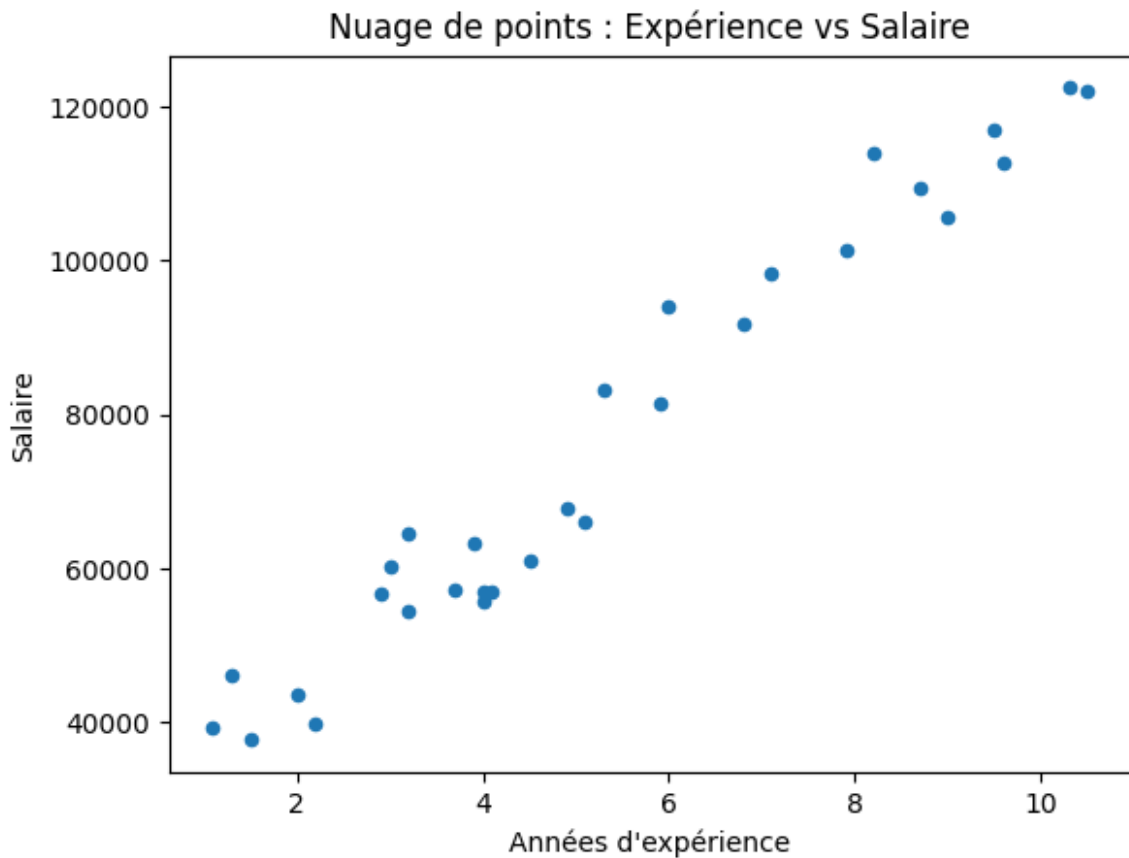


FIGURE 3 – Nuage des points pour le dataset Experience-Salaire

3.4 Nuages des points du deuxième data set selon les proprietes « Features »

Voici les visualisations demandees :

```
1 # Question 4 : Nuages des points pour le dataset assurance
2 # Selection des colonnes numeriques pour le scatter matrix
3 features = ['age', 'bmi', 'children', 'charges']
4 scatter_matrix(insurance[features], figsize=(10, 8), diagonal='hist', alpha
5               =0.7)
6
7 plt.suptitle('Scatter Matrix : Assurance', fontsize=16)
8 plt.show()
```

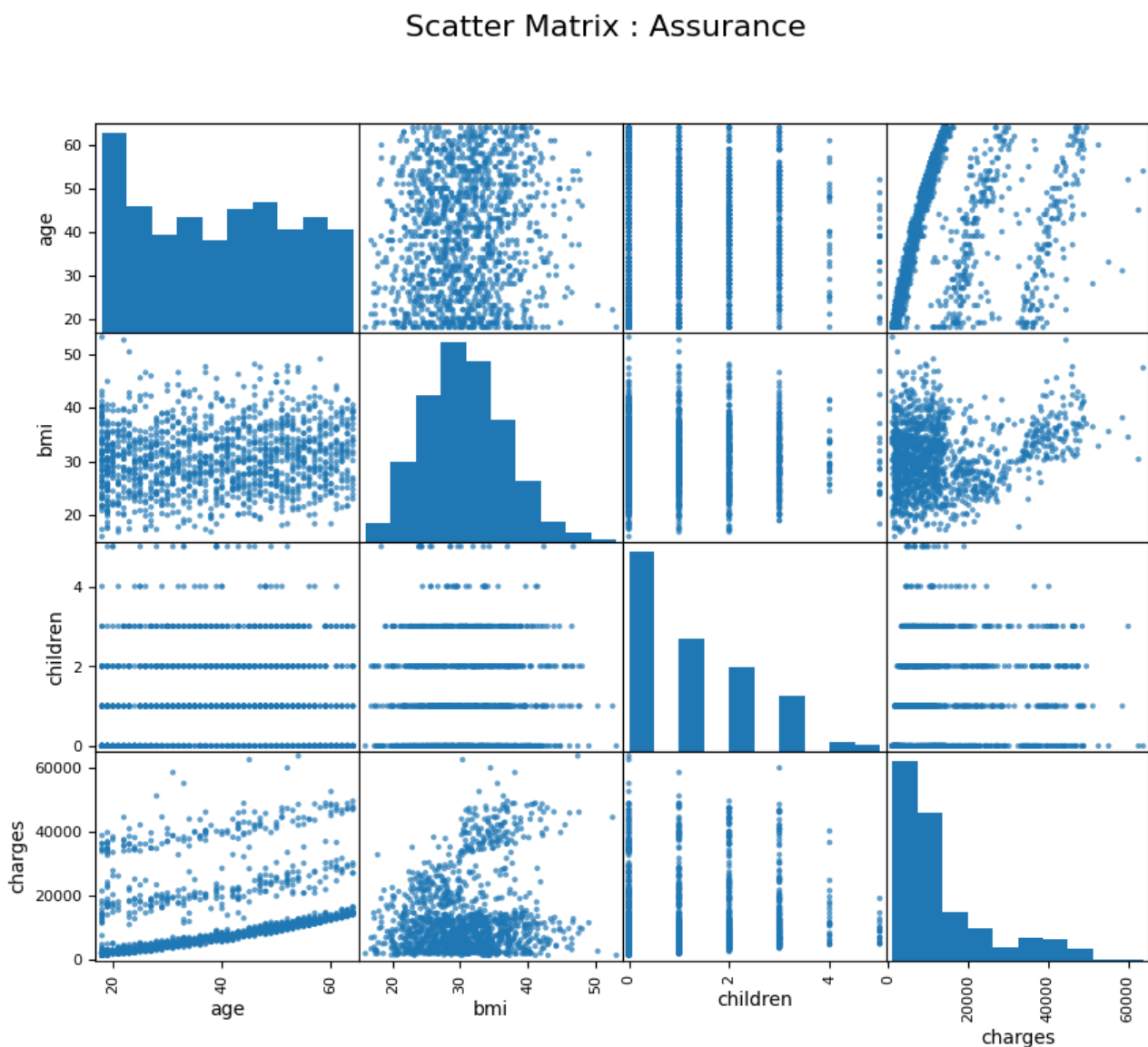


FIGURE 4 – Nuages des points du deuxième data set selon les proprietes « Features »

4 Partie 2 : Regression lineaire simple (Experience vs Salaire)

4.1 Entraînement de modèle

Pour predire les salaires en fonction des anneés d'experience, nous avons applique une regression lineaire simple :

```
1 # Partie 2 : Regression lineaire simple (Experience vs Salaire)
2 # Question 1 : Entra ner le mod le
3
4 from sklearn.model_selection import train_test_split
5 from sklearn.linear_model import LinearRegression
6
7 # Separation des donn es en variables explicatives (X) et cible (y)
8 X_exp = experience_salary[['YearsExperience']]
9 y_exp = experience_salary['Salary']
10
11 # Division des donn es en ensembles d'entra nement et de test
12 X_train, X_test, y_train, y_test = train_test_split(X_exp, y_exp, test_size
    =0.2, random_state=42)
13
14 # Creation et entra nement du mod le
15 model = LinearRegression()
16 model.fit(X_train, y_train)
17
18 # Affichage des coefficients du mod le
19 print(f"Coefficient de la regression : {model.coef_[0]}")
20 print(f"Intercept de la regression : {model.intercept_}")
```

Les visualisations et resultats d'evaluation sont inclus dans la suite du rapport.

— Coefficient de la regression : 9423.8153

— Intercept de la regression : 25321.5830

4.2 Prediction sur le dataset de test

```
1 # Question 2 : Prediction sur le dataset de test
2
3 y_pred = model.predict(X_test)
4
5 # Affichage des predictions
6 print("Predictions sur le dataset de test :")
7 print(y_pred)
```

— Predictions sur le dataset de test : [115790.21011287 71498.27809463 102596.86866063
75267.80422384 55477.79204548 60189.69970699]

4.3 Visualisation des resultats de la regression

```
1 # Question 3 : Visualisation des resultats de la regression
2 plt.scatter(X_test, y_test, color='blue', label='donn es reelles')
3 plt.plot(X_test, y_pred, color='red', label='Regression lineaire')
4 plt.title('Regression lineaire : Experience vs Salaire')
5 plt.xlabel('Annees d\'experience')
6 plt.ylabel('Salaire')
7 plt.legend()
8 plt.show()
```

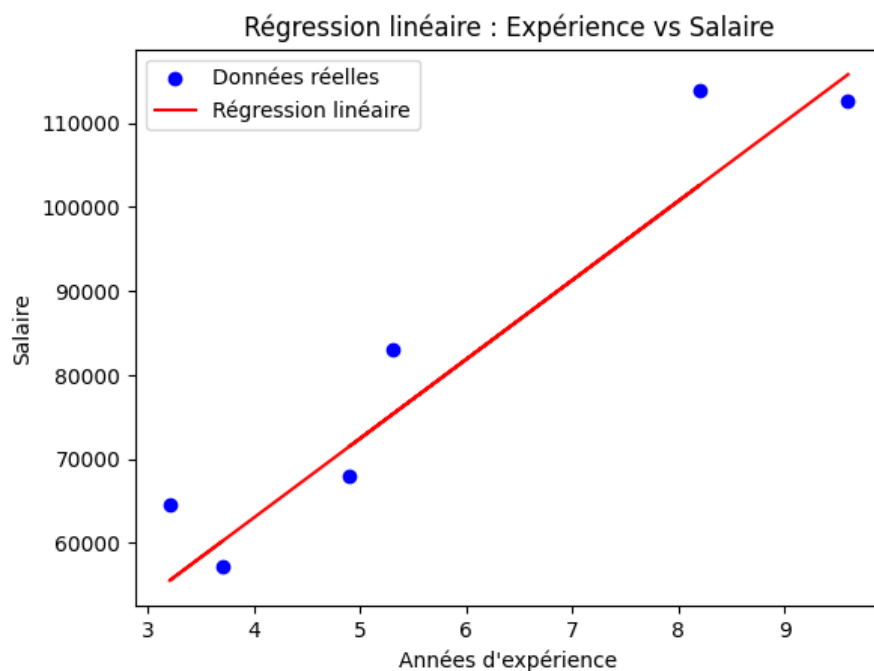


FIGURE 5 – Visualisation des resultats de la regression

4.4 Evaluation du model

```
1 # Question 4 : Evaluation du model
2 from sklearn.metrics import mean_squared_error, mean_absolute_error
3
4 # Calcul des metriques
5 mse = mean_squared_error(y_test, y_pred)
6 rmse = np.sqrt(mse)
7 mae = mean_absolute_error(y_test, y_pred)
8
9 # Affichage des resultats
10 print("----- evaluation du mod le : -----")
11 print(f"Mean Squared Error (MSE) : {mse}")
12 print(f"Root Mean Squared Error (RMSE) : {rmse}")
13 print(f"Mean Absolute Error (MAE) : {mae}")
```

```
----- Évaluation du modèle : -----
Mean Squared Error (MSE) : 49830096.85590839
Root Mean Squared Error (RMSE) : 7059.04362190151
Mean Absolute Error (MAE) : 6286.453830757749
```

FIGURE 6 – Evaluation du model

4.4.1 Interpretation

Les metriques calculees pour evaluer les performances du modèle sont les suivantes :

- **Mean Squared Error (MSE)** : 49,830,096.86
- **Root Mean Squared Error (RMSE)** : 7,059.04
- **Mean Absolute Error (MAE)** : 6,286.45

Ces resultats indiquent un ecart notable entre les predictions du modèle et les valeurs reelles. Une optimisation ou des ajustements supplementaires sont necessaires pour ameliorer la performance du modèle.

5 Partie 3 : Regression multiple (Assurance)

5.1 EDA (Exploratory Data Analysis)

EDA (Exploratory Data Analysis) :

```
1 # Question 1 : EDA (Exploratory Data Analysis)
2
3 # Aperçu des données
4 print("----- Aperçu des données : -----")
5 print(insurance.head())
6
7 # Informations générales sur les données
8 print("\n----- Informations générales : -----")
9 print(insurance.info())
```

```
----- Aperçu des données : -----
   age  sex  bmi  children  smoker  region  charges
0   19 female  27.900         0    yes southwest  16884.92400
1   18  male  33.770         1     no  southeast  1725.55230
2   28  male  33.000         3     no  southeast  4449.46200
3   33  male  22.705         0     no northwest  21984.47061
4   32  male  28.880         0     no northwest   3866.85520
```

```
----- Informations générales : -----
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1338 entries, 0 to 1337
Data columns (total 7 columns):
#   Column      Non-Null Count  Dtype
---  ---
0   age         1338 non-null   int64
1   sex         1338 non-null   object
2   bmi         1338 non-null   float64
3   children    1338 non-null   int64
4   smoker      1338 non-null   object
5   region      1338 non-null   object
6   charges     1338 non-null   float64
dtypes: float64(2), int64(2), object(3)
memory usage: 73.3+ KB
None
```

```
1 # Verification des valeurs nulles
2 print("\n----- Nombre de valeurs nulles par colonne : -----")
3 print(insurance.isnull().sum())
4
5 # Statistiques descriptives
6 print("\n----- Statistiques descriptives : -----")
7 print(insurance.describe())
```

```

----- Nombre de valeurs nulles par colonne : -----
age      0
sex      0
bmi      0
children 0
smoker   0
region   0
charges  0
dtype: int64

----- Statistiques descriptives : -----

```

	age	bmi	children	charges
count	1338.000000	1338.000000	1338.000000	1338.000000
mean	39.207025	30.663397	1.094918	13270.422265
std	14.049960	6.098187	1.205493	12110.011237
min	18.000000	15.960000	0.000000	1121.873900
25%	27.000000	26.296250	0.000000	4740.287150
50%	39.000000	30.400000	1.000000	9382.033000
75%	51.000000	34.693750	2.000000	16639.912515
max	64.000000	53.130000	5.000000	63770.428010

```

1 # Analyse des variables categoriques
2 print("\n----- Distribution des colonnes categoriques : -----")
3 print(insurance['sex'].value_counts())
4 print(insurance['smoker'].value_counts())
5 print(insurance['region'].value_counts())

```

```

----- Distribution des colonnes catégoriques : -----
sex
male      676
female    662
Name: count, dtype: int64
smoker
no        1064
yes        274
Name: count, dtype: int64
region
southeast    364
southwest    325
northwest    325
northeast    324
Name: count, dtype: int64

```

```

1 # Encoder les colonnes categoriques en variables numeriques
2 insurance_encoded = pd.get_dummies(insurance, drop_first=True)
3
4 # Matrice de correlation pour les variables numeriques
5 corr_matrix = insurance_encoded.corr()
6 print("\n----- Matrice de correlation : -----")
7 print(corr_matrix)

```

```

8
9 # Visualisation : Heatmap de la matrice de correlation
10 sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', fmt=".2f")
11 plt.title("Heatmap de correlation (donn es encodees)")
12 plt.show()

```

```

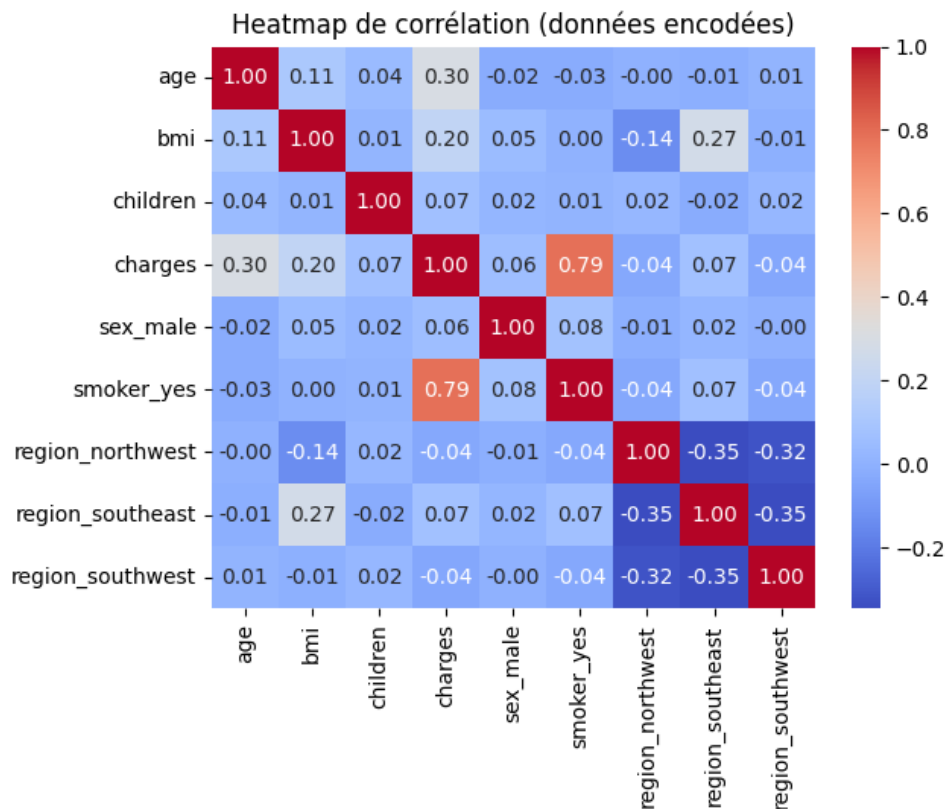
----- Matrice de corrélation : -----

```

	age	bmi	children	charges	sex_male \
age	1.000000	0.109272	0.042469	0.299008	-0.020856
bmi	0.109272	1.000000	0.012759	0.198341	0.046371
children	0.042469	0.012759	1.000000	0.067998	0.017163
charges	0.299008	0.198341	0.067998	1.000000	0.057292
sex_male	-0.020856	0.046371	0.017163	0.057292	1.000000
smoker_yes	-0.025019	0.003750	0.007673	0.787251	0.076185
region_northwest	-0.000407	-0.135996	0.024806	-0.039905	-0.011156
region_southeast	-0.011642	0.270025	-0.023066	0.073982	0.017117
region_southwest	0.010016	-0.006205	0.021914	-0.043210	-0.004184

	smoker_yes	region_northwest	region_southeast \
age	-0.025019	-0.000407	-0.011642
bmi	0.003750	-0.135996	0.270025
children	0.007673	0.024806	-0.023066
charges	0.787251	-0.039905	0.073982
sex_male	0.076185	-0.011156	0.017117
smoker_yes	1.000000	-0.036945	0.068498
region_northwest	-0.036945	1.000000	-0.346265
region_southeast	0.068498	-0.346265	1.000000
region_southwest	-0.036945	-0.320829	-0.346265

	region_southwest
age	0.010016
bmi	-0.006205
children	0.021914
charges	-0.043210
sex_male	-0.004184
smoker_yes	-0.036945
region_northwest	-0.320829
region_southeast	-0.346265
region_southwest	1.000000



5.2 Sélectionner 3 propriétés selon leurs degré d'importance

```

1 # Question 2 : Sélectionner 3 propriétés selon leurs degré d'importance
2
3
4 # 1. Calcul de la corrélation avec la variable cible
5 correlations = insurance_encoded.corr()
6 charges_corr = correlations['charges'].sort_values(ascending=False)
7 print("----- Corrélation avec 'charges' : -----")
8 print(charges_corr)
9
10 # 2. Affichage des 3 propriétés les plus corrélées
11 selected_features = charges_corr.index[1:4] # Exclut 'charges' lui-même
12 print("\n----- Propriétés sélectionnées selon la corrélation : -----",
13       selected_features.tolist())
14
15 # 3. Justification par Feature Importance avec un modèle
16 from sklearn.linear_model import LinearRegression
17 from sklearn.feature_selection import f_regression
18
19 X = insurance_encoded[selected_features]
20 y = insurance_encoded['charges']

```

```

20
21 # Entra nement du mod le
22 model = LinearRegression()
23 model.fit(X, y)
24
25 # Importance des variables
26 importances = pd.DataFrame({'Feature': selected_features, 'Importance':
    model.coef_})
27 print("\n----- Importance des variables selon la regression : -----")
28 print(importances)

```

```

----- Corrélation avec 'charges' : -----
charges          1.000000
smoker_yes       0.787251
age              0.299008
bmi              0.198341
region_southeast 0.073982
children         0.067998
sex_male         0.057292
region_northwest -0.039905
region_southwest -0.043210
Name: charges, dtype: float64

----- Propriétés sélectionnées selon la corrélation : ----- ['smoker_yes', 'age', 'bmi']

----- Importance des variables selon la régression : -----

```

	Feature	Importance
0	smoker_yes	23823.684495
1	age	259.547492
2	bmi	322.615133

5.3 Technique de standardisation

```

1 # Question 3 : technique de standardisation
2
3 from sklearn.preprocessing import StandardScaler
4
5 # Proprietes selectionnees
6 selected_features = ['age', 'bmi', 'smoker_yes']
7
8 # Standardisation des proprietes choisies
9 scaler = StandardScaler()
10 X_standardized = scaler.fit_transform(insurance_encoded[selected_features])
11
12 # Creation d'un DataFrame pour visualisation des resultats
13 X_standardized_df = pd.DataFrame(X_standardized, columns=selected_features)
14 print("\n----- Proprietes standardisees : -----")
15 print(X_standardized_df.head())

```



```

----- Propriétés standardisées : -----
      age      bmi  smoker_yes
0 -1.438764 -0.453320   1.970587
1 -1.509965  0.509621  -0.507463
2 -0.797954  0.383307  -0.507463
3 -0.441948 -1.305531  -0.507463
4 -0.513149 -0.292556  -0.507463

```

5.4 Entraîner le modèle

```

1 # Question 4 : Entraîner le modèle
2
3 from sklearn.model_selection import train_test_split
4 from sklearn.linear_model import LinearRegression
5
6 # Selection des variables explicatives et de la cible
7 X = X_standardized # Variables standardisées sélectionnées ('age', 'bmi', '
   smoker_yes')
8 y = insurance_encoded['charges']
9
10 # Division des données en ensembles d'entraînement et de test
11 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
   random_state=42)
12
13 # Creation et entraînement du modèle de regression lineaire
14 model = LinearRegression()
15 model.fit(X_train, y_train)
16
17 # Affichage des coefficients
18 print("\nCoefficients du modèle :")
19 for feature, coef in zip(['age', 'bmi', 'smoker_yes'], model.coef_):
20     print(f"{feature}: {coef:.4f}")
21
22 print(f"\nIntercept : {model.intercept_:.4f}")

```

```
Coefficients du modèle :  
age: 3643.3408  
bmi: 1990.0105  
smoker_yes: 9554.0342  
  
Intercept : 13321.2930
```

5.5 Prediction

```
1 # Question 5 : Predire  
2  
3 # Predire les donn es du dataset de test  
4 y_pred = model.predict(X_test)  
5  
6 # Affichage des predictions et des valeurs reelles  
7 predictions_df = pd.DataFrame({  
8     'Valeurs reelles': y_test.values,  
9     'Valeurs predites': y_pred  
10 })  
11  
12 print("\nPredictions sur le dataset de test :")  
13 print(predictions_df.head())
```

```
Prédictions sur le dataset de test :  
Valeurs réelles  Valeurs prédites  
0      9095.06825      8184.041468  
1      5272.17580      7431.001001  
2     29330.98315     37346.437099  
3      9301.89355      8629.528305  
4     33750.29180     27316.654810
```

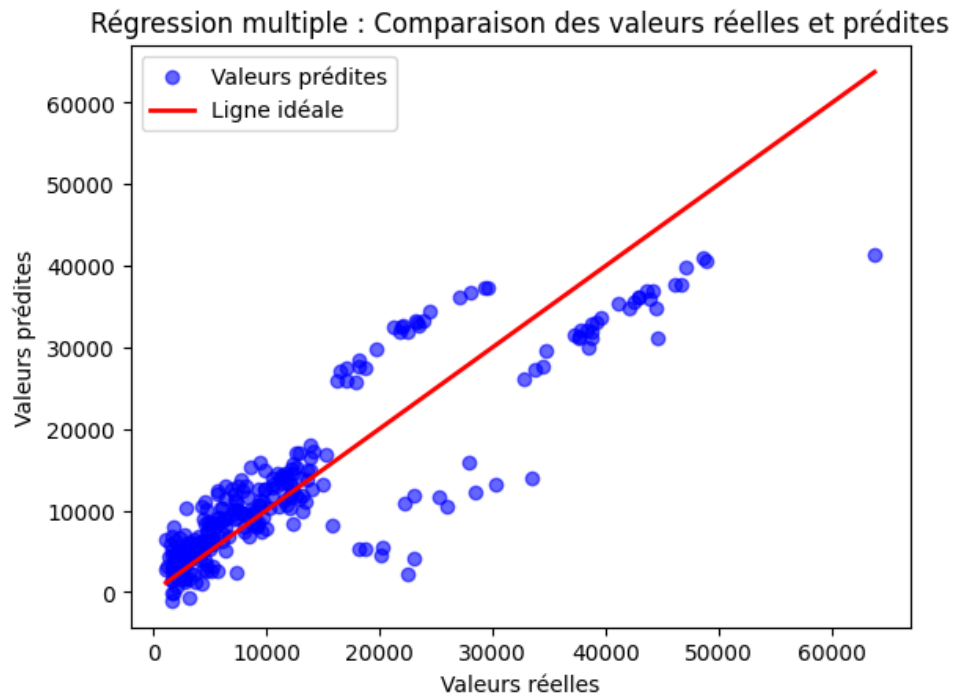
5.6 Visualisation

```
1 # Question 6 : Visualisation  
2  
3 # Visualisation des valeurs reelles vs valeurs predites  
4 plt.scatter(y_test, y_pred, alpha=0.6, color='blue', label='Valeurs predites  
    ,')
```

```

5 plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], color='
    red', lw=2, label='Ligne ideale')
6 plt.title('Regression multiple : Comparaison des valeurs reelles et predites
    ')
7 plt.xlabel('Valeurs reelles')
8 plt.ylabel('Valeurs predites')
9 plt.legend()
10 plt.show()

```



5.7 Evaluation de modèle en utilisant ces trois methodes

```

1 from sklearn.metrics import mean_squared_error, mean_absolute_error,
    r2_score
2 import numpy as np
3
4 # Calcul des metriques
5 mse = mean_squared_error(y_test, y_pred)
6 rmse = np.sqrt(mse)
7 mae = mean_absolute_error(y_test, y_pred)
8 r2 = r2_score(y_test, y_pred)
9
10 # Affichage des resultats
11 print("\n----- evaluation du mod le : -----")
12 print(f"Mean Squared Error (MSE) : {mse:.2f}")

```

```
----- Évaluation du modèle : -----  
Mean Squared Error (MSE) : 34512843.88  
Root Mean Squared Error (RMSE) : 5874.76  
Mean Absolute Error (MAE) : 4260.56  
Coefficient de détermination ( $R^2$ ) : 0.78
```

```
13 print(f"Root Mean Squared Error (RMSE) : {rmse:.2f}")  
14 print(f"Mean Absolute Error (MAE) : {mae:.2f}")  
15 print(f"Coefficient de determination (R ) : {r2:.2f}")
```

5.7.1 Interpretation

Les metriques calculees pour evaluer les performances du modèle sont les suivantes :

- **Mean Squared Error (MSE)** : 34,512,843.88
- **Root Mean Squared Error (RMSE)** : 5,874.76
- **Mean Absolute Error (MAE)** : 4,260.56
- **Coefficient de determination (R^2)** : 0.78

Analyse des resultats : Le coefficient de determination ($R^2 = 0.78$) indique que 78% de la variance des données est expliquée par le modèle, ce qui reflète une bonne capacité prédictive. Les erreurs (MSE, RMSE, MAE) sont raisonnablement faibles, mais une optimisation supplémentaire pourrait encore améliorer les prédictions.

6 Partie 4 : Regression polynomiale (China GDP)

6.1 Comparaison des modèles

Nous avons entraîné et comparé des modèles de regression lineaire et polynomiale pour predire l'évolution du PIB chinois :

```
1 # Question 1 : Entraîner les modèles
2 from sklearn.linear_model import LinearRegression
3 from sklearn.preprocessing import PolynomialFeatures
4 from sklearn.model_selection import train_test_split
5
6 # Separation des variables explicatives (X) et cible (y)
7 X_gdp = china_gdp[['Year']]
8 y_gdp = china_gdp['Value']
9
10 # Division des données en ensembles d'entraînement et de test
11 X_train, X_test, y_train, y_test = train_test_split(X_gdp, y_gdp, test_size
    =0.2, random_state=42)
12
13 # Modèle de regression lineaire
14 linear_model = LinearRegression()
15 linear_model.fit(X_train, y_train)
16
17 # Modèle de regression polynomiale (degré 4)
18 degree = 4
19 poly_features = PolynomialFeatures(degree=degree)
20 X_train_poly = poly_features.fit_transform(X_train)
21 X_test_poly = poly_features.transform(X_test)
22
23 poly_model = LinearRegression()
24 poly_model.fit(X_train_poly, y_train)
25
26 # Affichage des coefficients pour chaque modèle
27 print("\nCoefficients de la regression lineaire :")
28 print(f"Coef : {linear_model.coef_}, Intercept : {linear_model.intercept_}")
29
30 print("\nCoefficients de la regression polynomiale :")
31 print(f"Coef : {poly_model.coef_}, Intercept : {poly_model.intercept_}")
```

```

Coefficients de la régression linéaire :
Coef : [1.09727972e+11], Intercept : -216579790383883.8

Coefficients de la régression polynomiale :
Coef : [ 0.00000000e+00  6.04170555e+08  8.00117307e+11 -5.38980266e+08
 1.02111940e+05], Intercept : -5.224179684388893e+17

```

6.2 Prediction des données d'un data set de test pour les deux modèles

```

1 # Question 2 : Prediction des données d'un data set de test pour les deux
  modèles .
2
3 # Prediction avec le modèle de regression lineaire
4 linear_predictions = linear_model.predict(X_test)
5
6 # Prediction avec le modèle de regression polynomiale
7 poly_predictions = poly_model.predict(X_test_poly)
8
9 # Affichage des predictions
10 print("\n----- Predictions du modèle de regression lineaire : -----")
11 print(linear_predictions[:5])
12
13 print("\n----- Predictions du modèle de regression polynomiale : -----")
14 print(poly_predictions[:5])

```

```

----- Prédications du modèle de régression linéaire : -----
[ 1.88860122e+12 -9.64326039e+11  1.99832920e+12 -8.65022655e+10
 5.71865565e+11]

----- Prédications du modèle de régression polynomiale : -----
[-5.04822941e+10  2.41385689e+11  1.06164358e+10  5.04155164e+11
 2.24894391e+11]

```

6.3 Visualisation des resultats des deux modèles

```

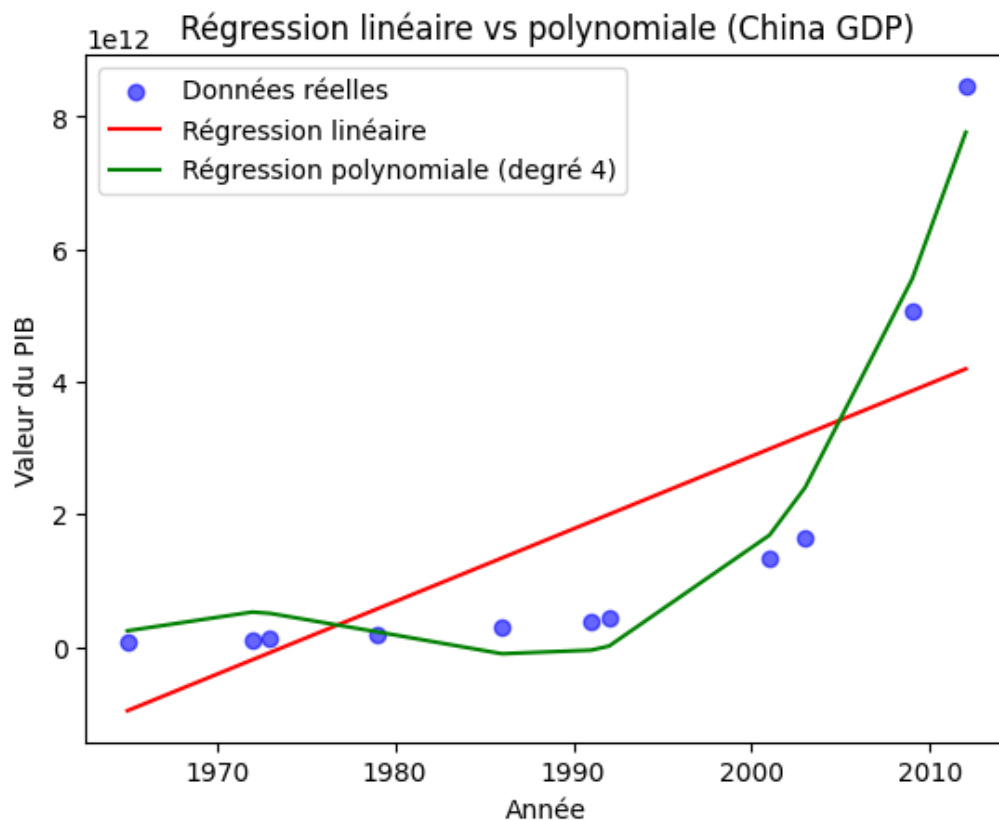
1 # Question 3 : Visualisation des resultats des deux modèles
2 # Trier les données de test et leurs predictions pour eviter des
  oscillations dans la visualisation
3 sorted_indices = np.argsort(X_test.values.flatten())
4 X_test_sorted = X_test.values[sorted_indices]
5 linear_predictions_sorted = linear_predictions[sorted_indices]
6 poly_predictions_sorted = poly_predictions[sorted_indices]

```

```

7
8 # Visualisation corrigee
9 plt.scatter(X_test, y_test, color='blue', label='donn es reelles', alpha
    =0.6)
10
11 # Tracer les lignes de prediction trieess
12 plt.plot(X_test_sorted, linear_predictions_sorted, color='red', label='
    Regression lineaire')
13 plt.plot(X_test_sorted, poly_predictions_sorted, color='green', label='
    Regression polynomiale (degre 4)')
14
15 # Configuration du graphique
16 plt.title('Regression lineaire vs polynomiale (China GDP)')
17 plt.xlabel('Annee')
18 plt.ylabel('Valeur du PIB')
19 plt.legend()
20 plt.show()

```



6.4 Evaluation des deux modèles en utilisant ces trois methodes

```
1 # Question 4 :  evaluation des deux mod les en utilisant ces trois methodes
2
3 from sklearn.metrics import r2_score
4
5 # evaluation du mod le de regression lineaire
6 linear_mse = mean_squared_error(y_test, linear_predictions)
7 linear_rmse = np.sqrt(linear_mse)
8 linear_mae = mean_absolute_error(y_test, linear_predictions)
9 linear_r2 = r2_score(y_test, linear_predictions)
10
11 # evaluation du mod le de regression polynomiale
12 poly_mse = mean_squared_error(y_test, poly_predictions)
13 poly_rmse = np.sqrt(poly_mse)
14 poly_mae = mean_absolute_error(y_test, poly_predictions)
15 poly_r2 = r2_score(y_test, poly_predictions)
16
17 # Affichage des resultats
18 print("\n----- evaluation du mod le de regression lineaire : -----")
19 print(f"Mean Squared Error (MSE) : {linear_mse:.2f}")
20 print(f"Root Mean Squared Error (RMSE) : {linear_rmse:.2f}")
21 print(f"Mean Absolute Error (MAE) : {linear_mae:.2f}")
22 print(f"Coefficient de determination (R ) : {linear_r2:.2f}")
23
24 print("\n----- evaluation du mod le de regression polynomiale : -----")
25 print(f"Mean Squared Error (MSE) : {poly_mse:.2f}")
26 print(f"Root Mean Squared Error (RMSE) : {poly_rmse:.2f}")
27 print(f"Mean Absolute Error (MAE) : {poly_mae:.2f}")
28 print(f"Coefficient de determination (R ) : {poly_r2:.2f}")
```

```
----- Évaluation du modèle de régression linéaire : -----
Mean Squared Error (MSE) : 2909722992049542559432704.00
Root Mean Squared Error (RMSE) : 1705791016522.70
Mean Absolute Error (MAE) : 1341446144991.78
Coefficient de détermination (R²) : 0.56
```

```
----- Évaluation du modèle de régression polynomiale : -----
Mean Squared Error (MSE) : 207534130192092768501760.00
Root Mean Squared Error (RMSE) : 455559140169.63
Mean Absolute Error (MAE) : 413445600606.50
Coefficient de détermination (R²) : 0.97
```


6.4.1 Interpretation

Les resultats de l'évaluation des deux modèles sont les suivants :

— **Modèle de regression lineaire :**

- **Mean Squared Error (MSE)** : 2,909,722,992,049,542,559,432,704.00
- **Root Mean Squared Error (RMSE)** : 1,705,791,016,522.70
- **Mean Absolute Error (MAE)** : 1,341,446,144,991.78
- **Coefficient de determination (R^2)** : 0.56

— **Modèle de regression polynomiale :**

- **Mean Squared Error (MSE)** : 207,534,130,192,092,768,501,760.00
- **Root Mean Squared Error (RMSE)** : 455,559,140,169.63
- **Mean Absolute Error (MAE)** : 413,445,600,606.50
- **Coefficient de determination (R^2)** : 0.97

Analyse des resultats : Le modèle de regression lineaire presente un $R^2 = 0.56$, ce qui signifie qu'il explique seulement 56% de la variance des données. Ses erreurs (MSE, RMSE, MAE) sont relativement elevees, indiquant une performance modeste.

En revanche, le modèle de regression polynomiale obtient un $R^2 = 0.97$, ce qui montre qu'il explique 97% de la variance des données, ce qui est un excellent resultat. Ses erreurs sont bien plus faibles que celles du modèle lineaire, ce qui suggère qu'il offre une meilleure prediction.

7 Conclusion

Cet atelier a permis d'explorer plusieurs types de regression et d'évaluer leur performance sur différents jeux de données. Les résultats obtenus montrent l'importance de choisir le modèle approprié pour les données disponibles.

Dans le cadre de cette étude, nous avons comparé deux modèles de regression : la regression lineaire et la regression polynomiale, en utilisant des metriques de performance telles que le Mean Squared Error (MSE), le Root Mean Squared Error (RMSE), le Mean Absolute Error (MAE) et le Coefficient de determination R^2 . Ces metriques nous ont permis de quantifier la precision et la qualite des predictions des modèles.

Le modèle de regression lineaire a montre une capacite d'ajustement relativement faible, avec un $R^2 = 0.56$, ce qui indique que plus de 40% de la variance des données reste inexpliquee. De plus, les erreurs (MSE, RMSE, MAE) elevees suggèrent que ce modèle n'est pas suffisamment adapte aux données analysees, ce qui peut être attribue à sa simplicite et à sa tendance à sous-estimer les relations complexes entre les variables.

En revanche, la regression polynomiale a demontre une performance bien superieure, avec un $R^2 = 0.97$, indiquant que presque 97% de la variance des données a ete expliquee par ce modèle. Les erreurs plus faibles (MSE, RMSE, MAE) confirment que le modèle polynomiale est capable de capturer de manière plus precise les relations non lineaires et complexes dans les données. Cela demontre l'importance de la flexibilite du modèle polynomial, bien que sa complexite puisse entraîner un surajustement (overfitting) dans certains cas.

Ces résultats renforcent l'idée que le choix du modèle depend non seulement de la nature des données, mais aussi du compromis entre la simplicite du modèle et sa capacite à capturer les relations sous-jacentes. Il est essentiel de tester plusieurs modèles et d'évaluer leur performance sur des metriques adaptees pour determiner la meilleure approche à adopter dans un contexte donne.

Ainsi, ce travail souligne l'importance de l'analyse prealable des données et du choix judicieux du modèle afin d'obtenir des résultats fiables et pertinents.

8 References

- Becoming Human - Linear Regression in Python
- Kaggle - Tutorial on Multiple Regression
- StackAbuse - Multiple Linear Regression