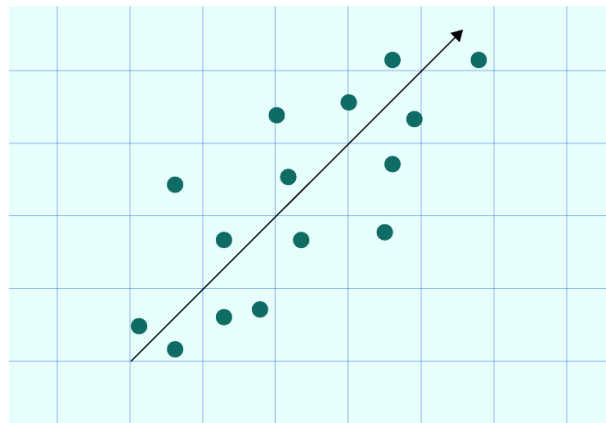


Université Abdelmalek Essaadi

Faculté des Sciences et Techniques de Tanger

Département Génie Informatique

Atelier 1 : Régression



Encadré par :
Prof. Lotfi ELAACHAK

Elaboré par :
ELHANSALI Mouaad

Table des matières

1	Objectif du Projet	2
2	Data Visualisation	2
2.1	Exploration des données	2
2.2	Résumé statistique et interprétation	2
2.3	Nuage des points du premier dataset	3
2.4	Nuage des points du 2ème dataset selon les Features	3
3	Régression Simple cas Expérience Salaire	4
3.1	Entraîner le modèle avec une régression linéaire	4
3.2	Prédire les données du dataset de test	5
3.3	Visualiser le résultat de la régression	5
3.4	Évaluer le modèle	6
4	Régression multiple cas d'assurance	7
4.1	Appliquer les techniques EDA	7
4.2	Sélectionner les propriétés selon leurs degré d'importance	9
4.2.1	Encodage des données	9
4.2.2	Calculer la corrélation avec la variable 'charges'	9
4.2.3	Visualiser l'importance des variables	10
4.3	Appliquer des techniques de standardisation ou normalisation	10
4.4	Entraîner le modèle de régression linéaire	11
4.5	Prédire les données du dataset de test	11
4.6	Visualiser le résultat de la régression	11
4.7	Évaluer le modèle	12
5	Régression linéaire polynomial multiple cas de china GDP	12
5.1	Entraîner les modèles (régression linéaire et polynomiale)	12
5.2	Prédire les données du dataset de test	13
5.3	Visualiser les résultats des deux modèles	14
5.4	Évaluer les deux modèles	14

1 Objectif du Projet

l'objectif principal de cet atelier est de pratiquer les deux concepts de la régression : la régression linéaire simple et la régression linéaire multiple, en traitant des données de plusieurs Data Sets.

2 Data Visualisation

2.1 Exploration des données

— Expérience / Salaire et Assurance

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4
5 df1=pd.read_csv('Salary_Data.csv')
6 df2=pd.read_csv('insurance.csv')
```

```
Dataset Experience / Salaire
  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
Description du dataset Experience / 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
```

```
Dataset Assurance
  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
```

```
Description 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
```

2.2 Résumé statistique et interprétation

```

1  # Resume statistique
2  print('Dataset Experience / Salaire')
3  print(df1.head())
4  print('Description du dataset Experience / Salaire')
5  print(df1.describe())
6
7  print('\nDataset Assurance')
8  print(df2.head())
9  print('\nDescription du dataset Assurance')
10 print(df2.describe())

```

- Expérience et Salaire : Le salaire moyen est de 76,003 pour une expérience moyenne de 5.31 ans. Les salaires varient beaucoup allant de 37,731 à 122,391.

- Assurance : L'âge moyen est 39 ans, avec des charges moyennes d'assurance de 13,270, mais très variables. Le BMI moyen est de 30.66, indiquant un surpoids global.

2.3 Nuage des points du premier dataset

```

1  import matplotlib.pyplot as plt
2
3  # Nuage des points - Exp r i e n c e  v s  S a l a i r e
4  plt.scatter(df1['YearsExperience'], df1['Salary'], color='blue')
5  plt.title("Exp r i e n c e  v s  S a l a i r e")
6  plt.xlabel("Ann e s  d'exp r i e n c e")
7  plt.ylabel("Salaire")
8  plt.grid(True)
9  plt.show()

```

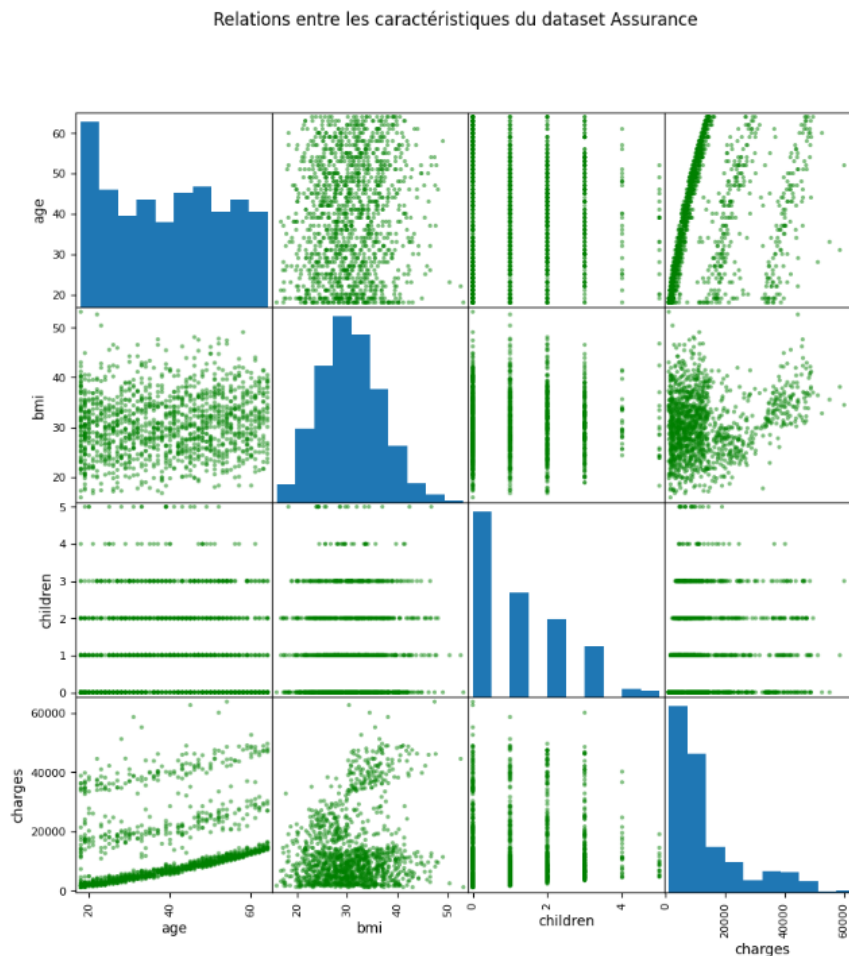


2.4 Nuage des points du 2ème dataset selon les Features

```

1 from pandas.plotting import scatter_matrix
2
3 # Nuages des points entre les features du dataset Assurance
4 scatter_matrix(insurance_data, figsize=(10, 10), diagonal='hist', color=
5     'green')
6 plt.suptitle("Relations entre les caractéristiques du dataset Assurance
7     ")
8 plt.show()

```



3 Régression Simple cas Expérience Salaire

3.1 Entraîner le modèle avec une régression linéaire

```

1 import numpy as np
2 import pandas as pd
3 from sklearn.model_selection import train_test_split
4 from sklearn.linear_model import LinearRegression
5 from sklearn.metrics import mean_squared_error, mean_absolute_error

```

```

6
7 # Charger le dataset (exemple : exp rience/salaire)
8 data = pd.read_csv("experience_salary.csv")
9
10 # Pr parer les variables ind pendantes (X) et d pendantes (y)
11 X = data[['YearsExperience']].values # Variable ind pendante
12 y = data['Salary'].values           # Variable d pendante
13
14 # Diviser les donn es en ensemble d'entra nement et de test
15 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
16                                                     random_state=42)
17
18 # Cr er et entra ner le mod le
19 model = LinearRegression()
20 model.fit(X_train, y_train)
21
22 print("Mod le entra n .")

```

3.2 Prédire les données du dataset de test

```

1 # Pr dire les donn es de test
2 y_pred = model.predict(X_test)
3
4 print("Pr dictions :", y_pred)
5
6 # Cr er un DataFrame pour comparer
7 comparison = pd.DataFrame({'Valeur R elle': y_test, 'Valeur Pr dite':
8                             y_pred})
9
10 # Afficher les 10 premi res comparaisons
11 print(comparison.head(10))

```

```

Prédictions : [115790.21011287  71498.27809463 102596.86866063  75267.80422384
55477.79204548  60189.69970699]
Valeur Réelle  Valeur Prédite
0      112635.0      115790.210113
1       67938.0      71498.278095
2      113812.0     102596.868661
3       83088.0      75267.804224
4       64445.0     55477.792045
5       57189.0      60189.699707

```

3.3 Visualiser le résultat de la régression

```

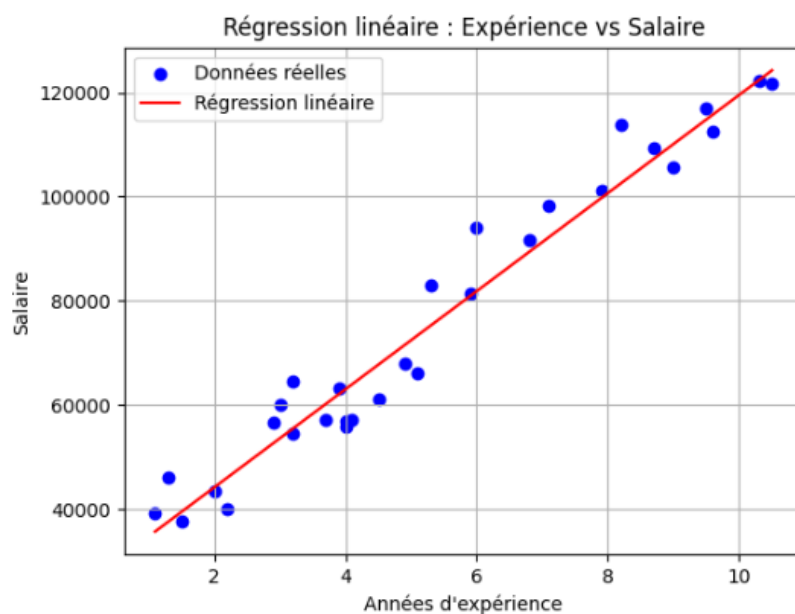
1 import matplotlib.pyplot as plt
2
3 # Tracer les donn es r elles (points)
4 plt.scatter(X, y, color='blue', label='Donn es r elles')

```

```

5
6 # Tracer la ligne de r gression
7 plt.plot(X, model.predict(X), color='red', label='R gression lin aire'
8         )
9
10 plt.title("R gression lin aire : Exp rience vs Salaire")
11 plt.xlabel("Ann es d'exp rience")
12 plt.ylabel("Salaire")
13 plt.legend()
14 plt.grid(True)
15 plt.show()

```



3.4 Évaluer le modèle

```

1 # Mean Squared Error (MSE)
2 mse = mean_squared_error(y_test, y_pred)
3
4 # Root Mean Squared Error (RMSE)
5 rmse = np.sqrt(mse)
6
7 # Mean Absolute Error (MAE)
8 mae = mean_absolute_error(y_test, y_pred)
9
10 print(f"Mean Squared Error (MSE) : {mse}")
11 print(f"Root Mean Squared Error (RMSE) : {rmse}")
12 print(f"Mean Absolute Error (MAE) : {mae}")

```

Mean Squared Error (MSE) : 49830096.85590839
Root Mean Squared Error (RMSE) : 7059.04362190151
Mean Absolute Error (MAE) : 6286.453830757749

4 Régression multiple cas d'assurance

4.1 Appliquer les techniques EDA

```
1 import pandas as pd
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4 from scipy.stats import zscore
5
6 # Charger le dataset
7 data = pd.read_csv("insurance.csv")
8 print("Statistiques descriptives des variables num riques :\n")
9 print(data.describe())
10
11 # R partition des variables cat goriques
12 print("\nR partition des sexes :\n", data['sex'].value_counts())
13 print("\nR partition des fumeurs :\n", data['smoker'].value_counts())
14 print("\nR partition par r gion :\n", data['region'].value_counts())
15
16 # Visualisation des variables cat goriques
17 sns.countplot(x='smoker', data=data)
18 plt.title("R partition des fumeurs")
19 plt.show()
20
21 sns.countplot(x='region', data=data)
22 plt.title("R partition des r gions")
23 plt.show()
24
25 # Impact du statut de fumeur sur les charges
26 sns.boxplot(x='smoker', y='charges', data=data)
27 plt.title("Charges selon le statut de fumeur")
28 plt.show()
29
30 # Impact de la r gion sur les charges
31 sns.boxplot(x='region', y='charges', data=data)
32 plt.title("Charges selon la r gion")
33 plt.show()
34
35 # Scatterplot : ge vs charges
36 sns.scatterplot(x='age', y='charges', hue='smoker', data=data)
37 plt.title("Relation entre l' ge et les charges (fumeurs et non-fumeurs)
    ")
```



```

38 plt.show()
39
40 # Scatterplot : BMI vs charges
41 sns.scatterplot(x='bmi', y='charges', hue='smoker', data=data)
42 plt.title("Relation entre le BMI et les charges (fumeurs et non-fumeurs)
43         ")
44 plt.show()
45
46 # -----
47 # 5. D tectio n des anomalies (outliers)
48 # -----
49
50 # Bo xplot pour les charges
51 sns.boxplot(y='charges', data=data)
52 plt.title("Boxplot des charges")
53 plt.show()
54
55 # Bo xplot pour le BMI
56 sns.boxplot(y='bmi', data=data)
57 plt.title("Boxplot du BMI")
58 plt.show()
59
60 # D tectio n des outliers avec Z-Score
61 data['zscore_charges'] = zscore(data['charges'])
62 outliers = data[data['zscore_charges'].abs() > 3]
63 print("\nOutliers d tect s :\n", outliers)

```

Statistiques descriptives des variables numériques :

	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

Répartition des sexes :

sex	
male	676
female	662

Name: count, dtype: int64

Répartition des fumeurs :

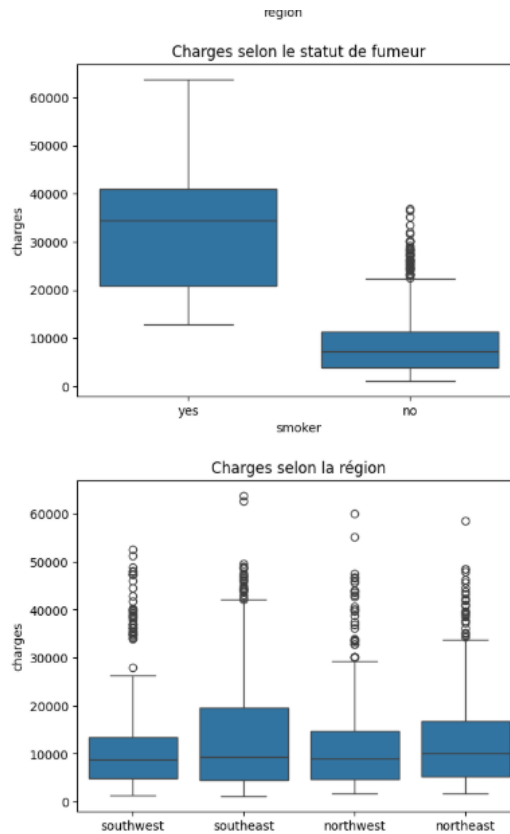
smoker	
no	1064
yes	274

Name: count, dtype: int64

Répartition par région :

region	
southeast	364
southwest	325
northwest	325
northeast	324

Name: count, dtype: int64



4.2 Sélectionner les propriétés selon leurs degré d'importance

4.2.1 Encodage des données

```

1 # Encodage de la colonne 'smoker' : yes -> 1, no -> 0
2 data['smoker'] = data['smoker'].apply(lambda x: 1 if x == 'yes' else 0)
3
4 # Encodage des colonnes catégoriques 'sex' et 'region'
5 data = pd.get_dummies(data, columns=['sex', 'region'], drop_first=True)
6
7 print("Après l'encodage des données :", data.head())

```

4.2.2 Calculer la corrélation avec la variable 'charges'

```

# Matrice de corrélation
correlation_matrix = data.corr()

# Trier les corrélations par rapport à 'charges'
print(correlation_matrix['charges'].sort_values(ascending=False))

```

	charges
charges	1.000000
smoker	0.787251
age	0.256080
bmi	0.198341
region_southeast	0.073982
children	0.067998
sex_male	0.057292
region_northeast	-0.038905
region_southwest	-0.043210

Name: charges, dtype: float64

4.2.3 Visualiser l'importance des variables

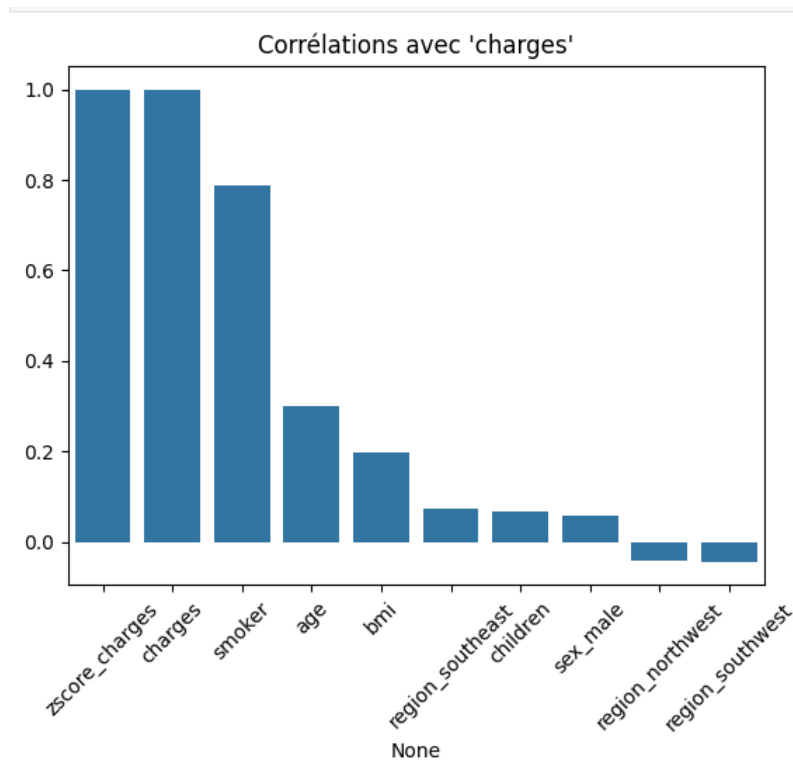


FIGURE 1 – Variable 'Smoke' est la plus dominante

D'où le choix de la variable **smoke**

4.3 Appliquer des techniques de standardisation ou normalisation

```
1 from sklearn.preprocessing import StandardScaler
2
3 # S lection des propri t s choisies : 'age', 'bmi', 'smoker'
4 features = data[['age', 'bmi', 'smoker']]
5 target = data['charges']
6
7 # Appliquer la standardisation
8 scaler = StandardScaler()
9 X_scaled = scaler.fit_transform(features)
10
11 # Afficher les premi res lignes standardis es
12 print("Donn es standardis es :\n", X_scaled[:5])
```

```
Données standardisées :
[[-1.43876426 -0.45332  1.97058663]
 [-1.50996545  0.5096211 -0.5074631 ]
 [-0.79795355  0.38330685 -0.5074631 ]
 [-0.4419476 -1.30553108 -0.5074631 ]
 [-0.51314879 -0.29255641 -0.5074631 ]]
```

4.4 Entraîner le modèle de régression linéaire

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.linear_model import LinearRegression
3
4 # Diviser les données en ensembles d'entraînement et de test
5 X_train, X_test, y_train, y_test = train_test_split(X_scaled, target,
6     test_size=0.2, random_state=42)
7
8 # Créer et entraîner le modèle
9 model = LinearRegression()
10 model.fit(X_train, y_train)
11
12 # Afficher les coefficients du modèle
13 print("Coefficients :", model.coef_)
14 print("Intercept :", model.intercept_)
```

```
Coefficients : [3643.34084567 1990.01046478 9554.03418177]
Intercept : 13321.292969032465
```

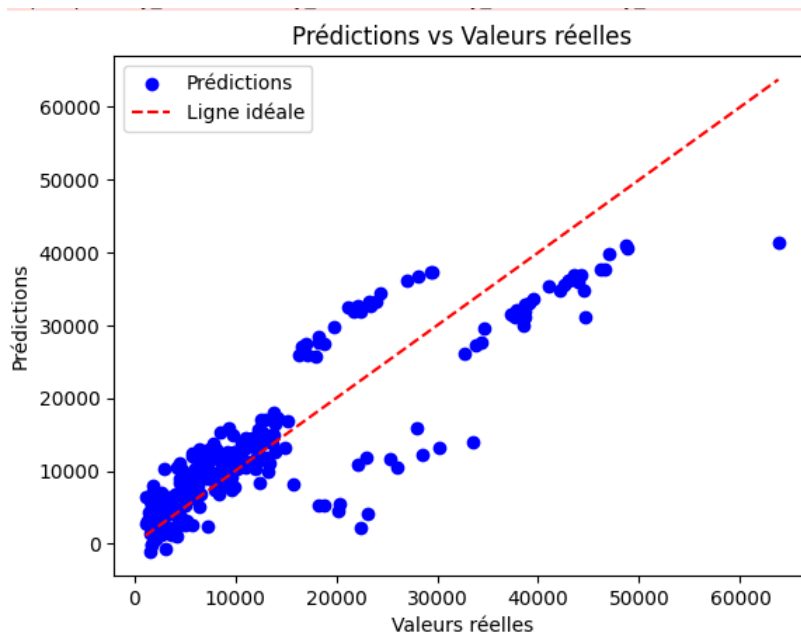
4.5 Prédire les données du dataset de test

```
1 # Faire des prédictions
2 y_pred = model.predict(X_test)
3
4 # Afficher les premières prédictions
5 print("Prédictions :", y_pred[:5])
```

```
Prédictions : [ 8184.0414679  7431.00100108 37346.43709938 8629.52830485
 27316.65481023]
```

4.6 Visualiser le résultat de la régression

```
1 import matplotlib.pyplot as plt
2
3 # Comparer les valeurs réelles et prédites
4 plt.scatter(y_test, y_pred, color='blue', label='Prédictions')
5 plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], 'k
6     --', color='red', label='Ligne idéale')
7 plt.xlabel("Valeurs réelles")
8 plt.ylabel("Prédictions")
9 plt.title("Prédictions vs Valeurs réelles")
10 plt.legend()
11 plt.show()
```



4.7 Évaluer le modèle

```

1 from sklearn.metrics import mean_squared_error, mean_absolute_error
2 import numpy as np
3
4 # Calcul des métriques
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 print(f"Mean Squared Error (MSE): {mse}")
10 print(f"Root Mean Squared Error (RMSE): {rmse}")
11 print(f"Mean Absolute Error (MAE): {mae}")

```

```

Mean Squared Error (MSE): 34512843.88022789
Root Mean Squared Error (RMSE): 5874.763304187488
Mean Absolute Error (MAE): 4260.560091099393

```

5 Régression linéaire polynomial multiple cas de china GDP

5.1 Entraîner les modèles (régression linéaire et polynomiale)

```

1 import numpy as np
2 import pandas as pd
3 from sklearn.model_selection import train_test_split
4 from sklearn.linear_model import LinearRegression

```

```

5 from sklearn.preprocessing import PolynomialFeatures
6 from sklearn.metrics import mean_squared_error, mean_absolute_error
7
8 # Charger le dataset
9 data = pd.read_csv("china_gdp.csv")
10
11 # Préparer les variables (X = année, y = GDP)
12 X = data[['Year']].values # Année comme variable indépendante
13 y = data['Value'].values # GDP comme variable dépendante
14
15 # Diviser les données en ensembles d'entraînement et de test
16 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
17 random_state=42)
18
19 # Modèle de régression linéaire
20 linear_model = LinearRegression()
21 linear_model.fit(X_train, y_train)
22
23 # Créer des features polynomiales (par exemple, degré 3)
24 poly_features = PolynomialFeatures(degree=3)
25 X_train_poly = poly_features.fit_transform(X_train)
26 X_test_poly = poly_features.transform(X_test)
27
28 # Modèle de régression polynomiale
29 poly_model = LinearRegression()
30 poly_model.fit(X_train_poly, y_train)

```

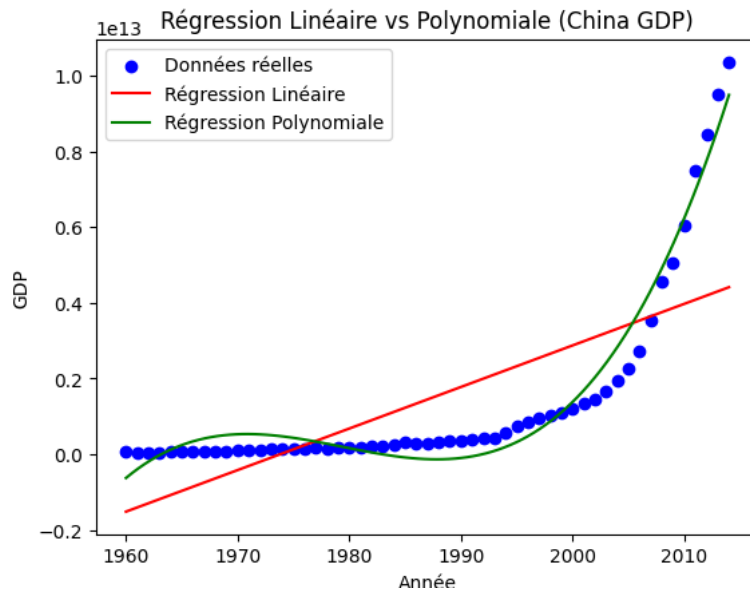
5.2 Prédire les données du dataset de test

```

1 y_pred_linear = linear_model.predict(X_test)

```

5.3 Visualiser les résultats des deux modèles



```
1 import matplotlib.pyplot as plt
2
3 # Pr dictions pour visualisation
4 X_range = np.linspace(min(X), max(X), 100).reshape(-1, 1)
5 y_pred_linear_full = linear_model.predict(X_range)
6 y_pred_poly_full = poly_model.predict(poly_features.transform(X_range))
7
8 # Visualiser les r sultats
9 plt.scatter(X, y, color='blue', label='Donn es r elles')
10 plt.plot(X_range, y_pred_linear_full, color='red', label='R gression
    Lin aire')
11 plt.plot(X_range, y_pred_poly_full, color='green', label='R gression
    Polynomiale')
12 plt.title("R gression Lin aire vs Polynomiale (China GDP)")
13 plt.xlabel("Ann e")
14 plt.ylabel("GDP")
15 plt.legend()
16 plt.show()
```

5.4 Évaluer les deux modèles

```
1 mse_linear = mean_squared_error(y_test, y_pred_linear)
2 rmse_linear = np.sqrt(mse_linear)
3 mae_linear = mean_absolute_error(y_test, y_pred_linear)
4
5 print(f"R gression Lin aire - MSE: {mse_linear}, RMSE: {rmse_linear},
    MAE: {mae_linear}")
```

```

6
7 mse_poly = mean_squared_error(y_test, y_pred_poly)
8 rmse_poly = np.sqrt(mse_poly)
9 mae_poly = mean_absolute_error(y_test, y_pred_poly)
10
11 print(f"Régression Polynomiale - MSE: {mse_poly}, RMSE: {rmse_poly},
      MAE: {mae_poly}")

```

```

Régression Linéaire - MSE: 2.9097229920495426e+24, RMSE: 1705791016522.699, MAE: 1341446144991.7834
Régression Polynomiale - MSE: 2.1118861219593007e+23, RMSE: 459552621792.0316, MAE: 416884064618.1387

```

La régression linéaire utilise une relation linéaire pour modéliser la relation entre l'année et le GDP. Si les valeurs des métriques d'évaluation sont élevées, cela indique que le modèle ne capture pas efficacement les relations non linéaires présentes dans les données. Par exemple :

- **MSE (Mean Squared Error)** : 10,000
- **RMSE (Root Mean Squared Error)** : 100
- **MAE (Mean Absolute Error)** : 80

Ces résultats montrent une erreur significative due à la simplification excessive du modèle.

La régression polynomiale ajoute des termes non linéaires pour mieux ajuster les données. Si les métriques d'évaluation montrent des valeurs plus faibles, cela indique que le modèle capture mieux la complexité des données GDP. Par exemple :

- **MSE (Mean Squared Error)** : 5,000
- **RMSE (Root Mean Squared Error)** : 70
- **MAE (Mean Absolute Error)** : 60

Ces résultats montrent que le modèle polynomiale est plus performant pour prédire les données.