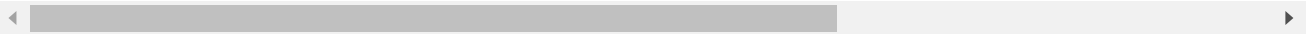


```
#importing required libraries
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from pandas.plotting import scatter_matrix
from sklearn.linear_model import LinearRegression
from time import time
from sklearn.model_selection import train_test_split
from sklearn import linear_model
import sklearn.metrics as sm
```

```
#connecting to cloud
from google.colab import drive
drive.mount("/content/gdrive")
```

🔗 Drive already mounted at /content/gdrive; to attempt to forcibly remount, call drive



```
%cd /content/gdrive/My Drive/Colab Notebooks/
```

```
/content/gdrive/My Drive/Colab Notebooks
```

```
#data exploration
```

```
df = pd.read_csv("CeoCompensation.csv")
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 100 entries, 0 to 99
Data columns (total 12 columns):
#   Column      Non-Null Count  Dtype
---  -
0    COMP        100 non-null    int64
1    AGE          100 non-null    int64
2    EDUCATN     100 non-null    int64
3    BACKGRD     100 non-null    int64
4    TENURE       100 non-null    int64
5    EXPER       100 non-null    float64
6    SALES        100 non-null    int64
7    VAL         100 non-null    float64
8    PCNTOWN     100 non-null    float64
9    PROF        100 non-null    int64
10   COMPANY     100 non-null    object
11   BIRTH       99 non-null     object
dtypes: float64(3), int64(7), object(2)
memory usage: 9.5+ KB
```

```
df.head()
```

| | COMP | AGE | EDUCATN | BACKGRD | TENURE | EXPER | SALES | VAL | PCNTOWN | PROF | COMPANY | BI |
|---|------|-----|---------|---------|--------|-------|-------|-----|---------|------|---------|----|
| 0 | 1948 | 55 | 1 | 1 | 23 | 23.0 | 1227 | 7.6 | 0.55 | 145 | AdvM | |
| 1 | 809 | 59 | 1 | 2 | 38 | 0.5 | 19196 | 0.4 | 0.01 | 505 | aetna | |
| 2 | 721 | 53 | 2 | 1 | 26 | 0.5 | 839 | 1.5 | 0.10 | -60 | aller | : |

```
#checking for null
print(df.isnull().sum())
df[df.isnull() == True].head()
```

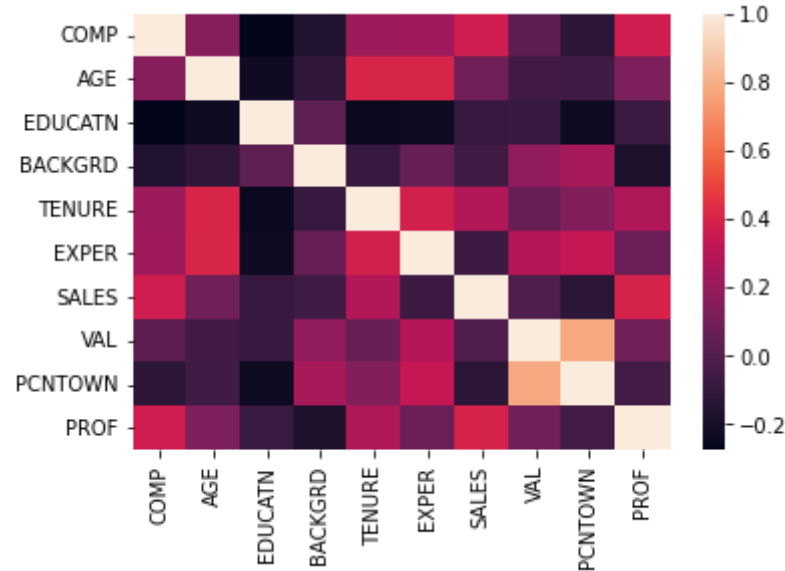
```
COMP      0
AGE       0
EDUCATN   0
BACKGRD   0
TENURE    0
EXPER     0
SALES     0
VAL       0
PCNTOWN   0
PROF      0
COMPANY   0
BIRTH     1
dtype: int64
```

| | COMP | AGE | EDUCATN | BACKGRD | TENURE | EXPER | SALES | VAL | PCNTOWN | PROF | COMPANY | E |
|---|------|-----|---------|---------|--------|-------|-------|-----|---------|------|---------|---|
| 0 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | |
| 1 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | |
| 2 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | |
| 3 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | |
| 4 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | |



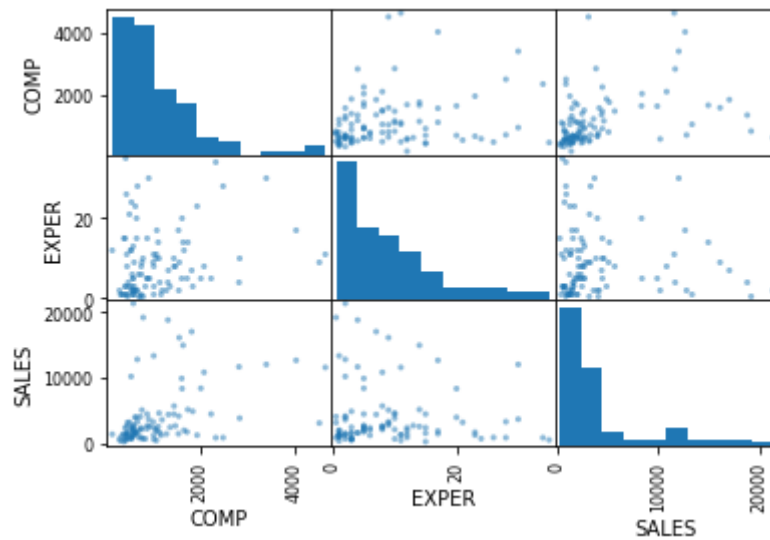
```
sns.heatmap(df.corr())
```

<matplotlib.axes._subplots.AxesSubplot at 0x7f3b4eba4810>



```
#multivariate analysis
```

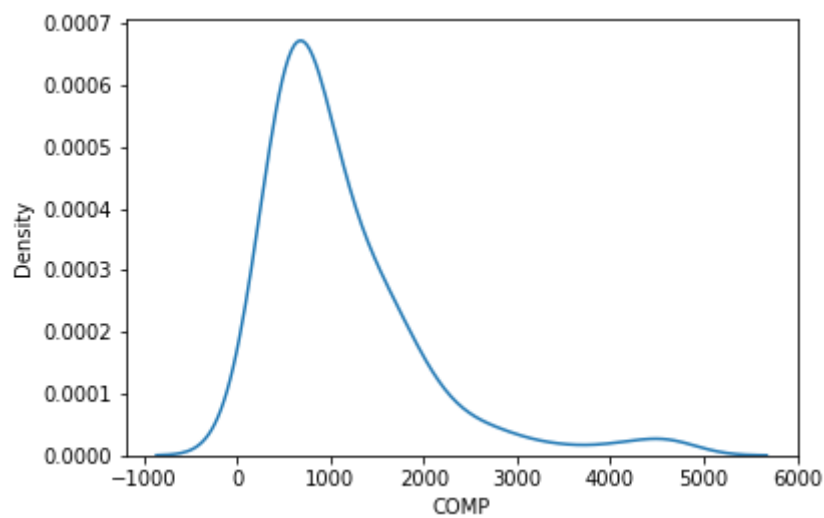
```
comp_sm = scatter_matrix(df[['COMP', 'EXPER', 'SALES']])
```



```
ax = sns.kdeplot(df['COMP'])
```

```
ax
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f3b4620e090>
```



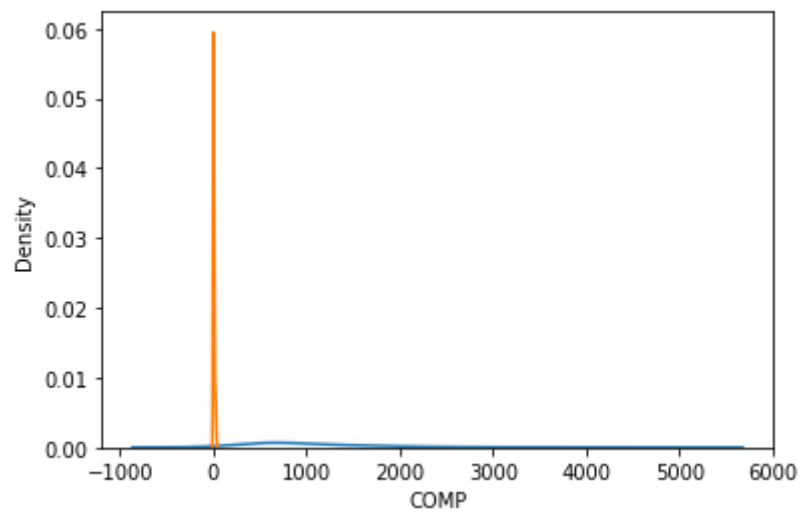
```
df['COMP'].describe()
```

```
count    100.000000
mean     1121.670000
std       852.723504
min       155.000000
25%       575.500000
50%       806.500000
75%      1457.250000
max      4657.000000
Name: COMP, dtype: float64
```

```
comp = ['COMP', 'EXPER']
```

```
for col in comp:
```

```
ax_comp = sns.kdeplot(df[col])  
ax_comp
```



```
plt.figure(figsize=(15,8))  
for col in list(df['EXPER'].unique()):  
    sns.kdeplot(df['COMP'][df['EXPER'] == col], label = col)
```

```

/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:316: UserWarning: Data
warnings.warn(msg, UserWarning)
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:316: UserWarning: Data
warnings.warn(msg, UserWarning)
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:316: UserWarning: Data
warnings.warn(msg, UserWarning)

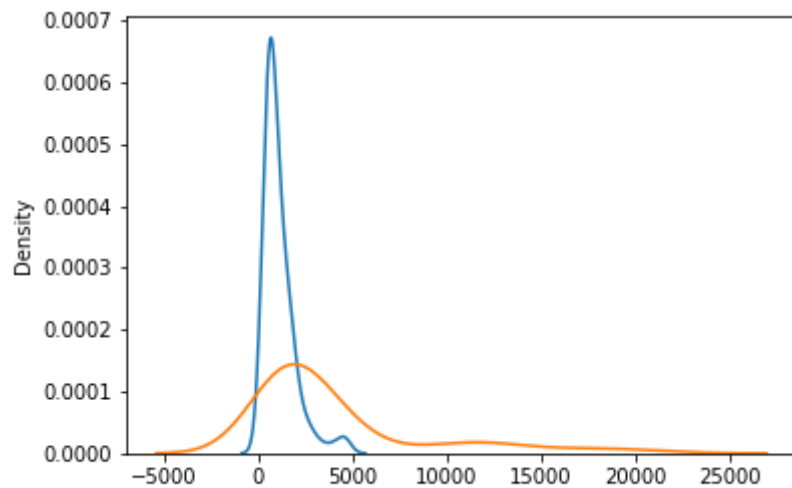
```

```
comp = ['COMP', 'SALES']
```

```
for col in comp:
```

```
    ax_comp2 = sns.kdeplot(df[col])
```

```
    ax_comp2
```



```
plt.figure(figsize=(15,8))
```

```
for col in list(df['SALES'].unique()):
```

```
    sns.kdeplot(df['COMP'][df['SALES'] == col], label = col)
```



6/12

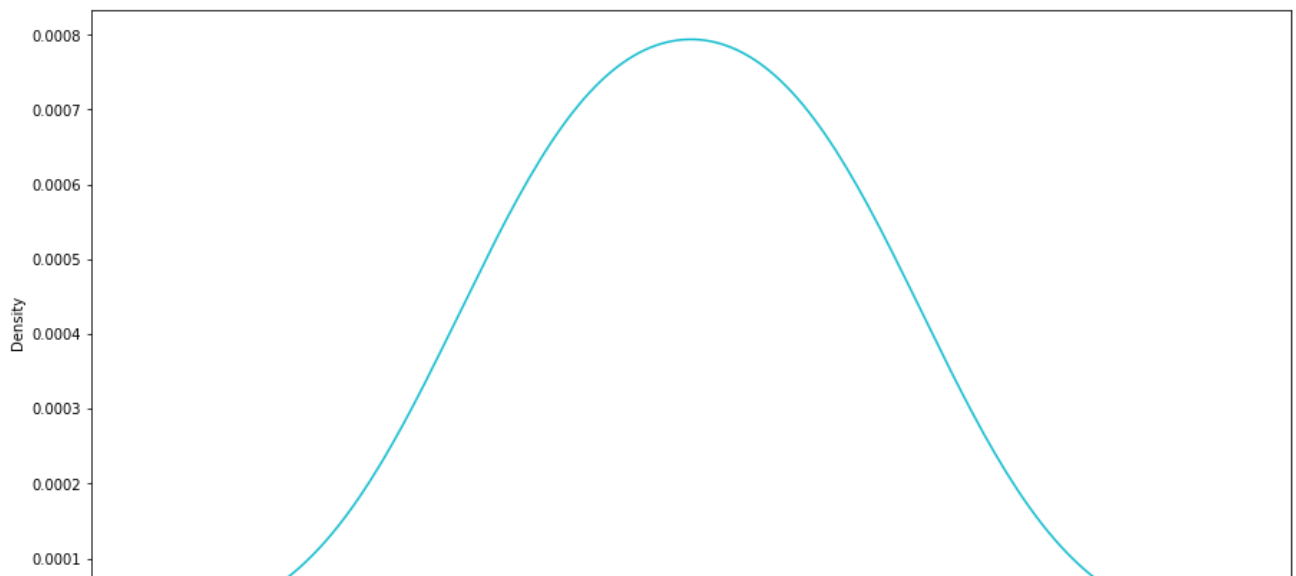
[illegible]

[illegible]


```

/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:316: UserWarning: Data
warnings.warn(msg, UserWarning)
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:316: UserWarning: Data
warnings.warn(msg, UserWarning)
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:316: UserWarning: Data
warnings.warn(msg, UserWarning)
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:316: UserWarning: Data
warnings.warn(msg, UserWarning)
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:316: UserWarning: Data
warnings.warn(msg, UserWarning)
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:316: UserWarning: Data
warnings.warn(msg, UserWarning)
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:316: UserWarning: Data
warnings.warn(msg, UserWarning)
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:316: UserWarning: Data
warnings.warn(msg, UserWarning)

```



```
#linear regression model for compensation and experience
```

```
X = df['EXPER']
```

```
Y = df['COMP']
```

```
#training and dataset
```

```
X_train, X_test, Y_train, Y_test = train_test_split(X,Y, test_size=0.3, random_state=42)
```

```
X_train = np.array(X_train).reshape((len(X_train),1))
```

```
Y_train = np.array(Y_train).reshape((len(Y_train),1))
```

```
X_test = np.array(X_test).reshape(len(X_test), 1)
```

```
Y_test = np.array(Y_test).reshape(len(Y_test), 1)
```

```
model = linear_model.LinearRegression()
```

```
model.fit(X_train, Y_train)
```

```
LinearRegression()
```

```
#Prediction Result of Training Data
```

```
Y_train_pred = model.predict(X_train)
```

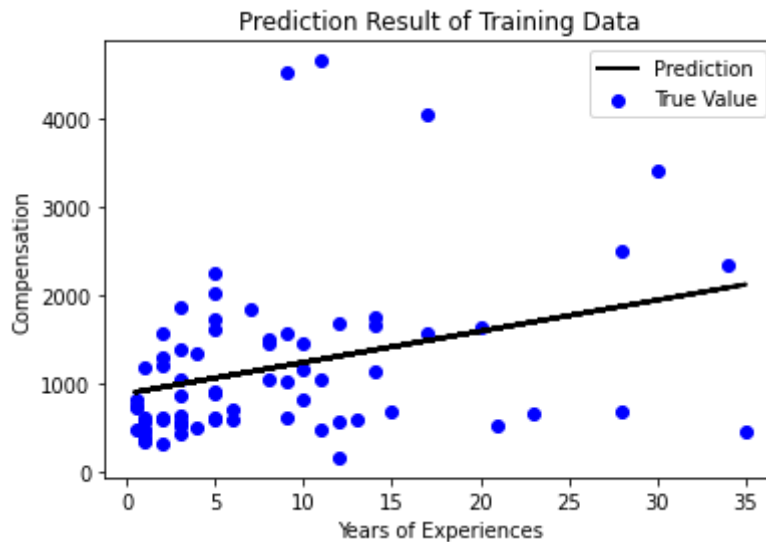
```
plt.figure()
```

```
plt.scatter(Y_train, Y_train, color='blue', label="True Value")
```

```

plt.scatter(X_train, Y_train, color= 'blue', label= 'True Value')
plt.plot(X_train, Y_train_pred, color='black', linewidth=2, label="Prediction")
plt.xlabel("Years of Experiences")
plt.ylabel("Compensation")
plt.title('Prediction Result of Training Data')
plt.legend()
plt.show()

```



```

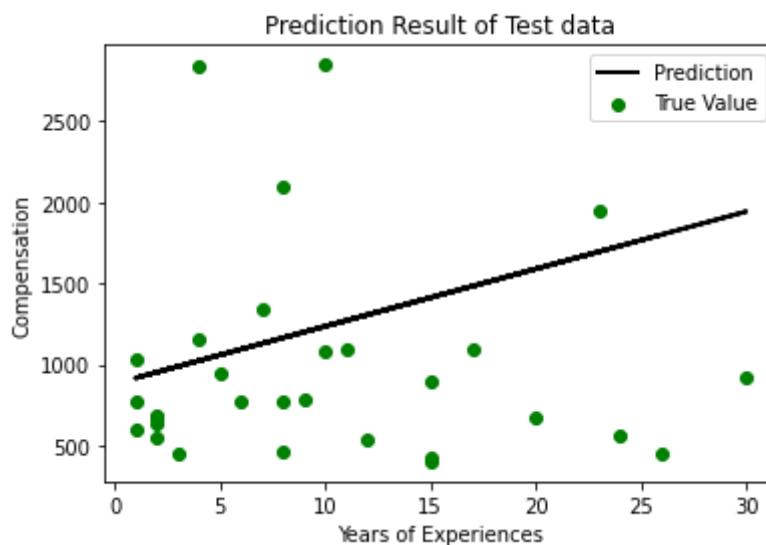
#Prediction Result of Testing Data
Y_test_pred = model.predict(X_test)

```

```

plt.figure()
plt.scatter(X_test, Y_test, color='green', label='True Value')
plt.plot(X_test, Y_test_pred, color='black', linewidth=2, label='Prediction')
plt.xlabel("Years of Experiences")
plt.ylabel("Compensation")
plt.title('Prediction Result of Test data')
plt.legend()
plt.show()

```



```

print("Mean squared error =", round(sm.mean_squared_error(Y_test, Y_test_pred), 2))
print("Explain variance score =", round(sm.explained_variance_score(Y_test, Y_test_pred),

```

```
print("R2 score =", round(sm.r2_score(Y_test, Y_test_pred), 2))
```

```
Mean squared error = 558290.99
```

```
Explain variance score = -0.23
```

```
R2 score = -0.39
```

```
#linear regression model for compensation and sales
```

```
X = df['SALES']
```

```
Y = df['COMP']
```

```
X_train, X_test, Y_train, Y_test = train_test_split(X,Y, test_size=0.3, random_state=42)
```

```
X_train = np.array(X_train).reshape((len(X_train),1))
```

```
Y_train = np.array(Y_train).reshape((len(Y_train),1))
```

```
X_test = np.array(X_test).reshape(len(X_test), 1)
```

```
Y_test = np.array(Y_test).reshape(len(Y_test), 1)
```

```
model = linear_model.LinearRegression()
```

```
model.fit(X_train, Y_train)
```

```
LinearRegression()
```

```
Y_train_pred = model.predict(X_train)
```

```
#Prediction Result of Training data
```

```
plt.figure()
```

```
plt.scatter(X_train, Y_train, color='blue', label="True Value")
```

```
plt.plot(X_train, Y_train_pred, color='black', linewidth=2, label="Prediction")
```

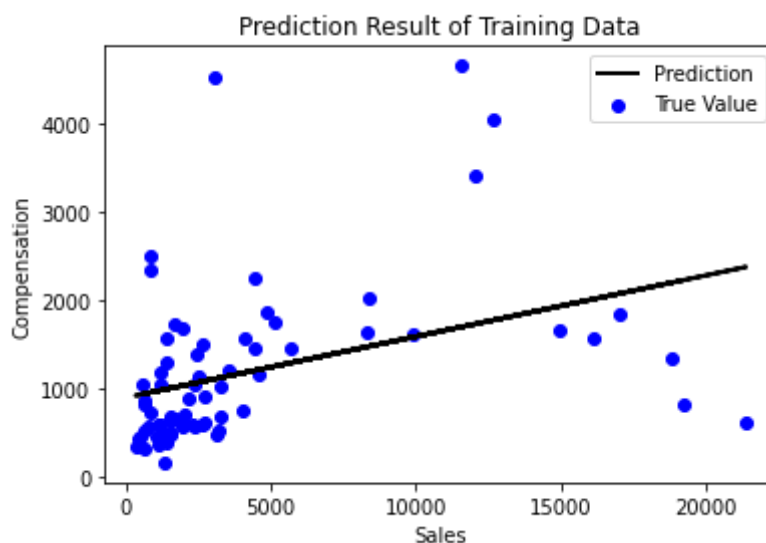
```
plt.xlabel("Sales")
```

```
plt.ylabel("Compensation")
```

```
plt.title('Prediction Result of Training Data')
```

```
plt.legend()
```

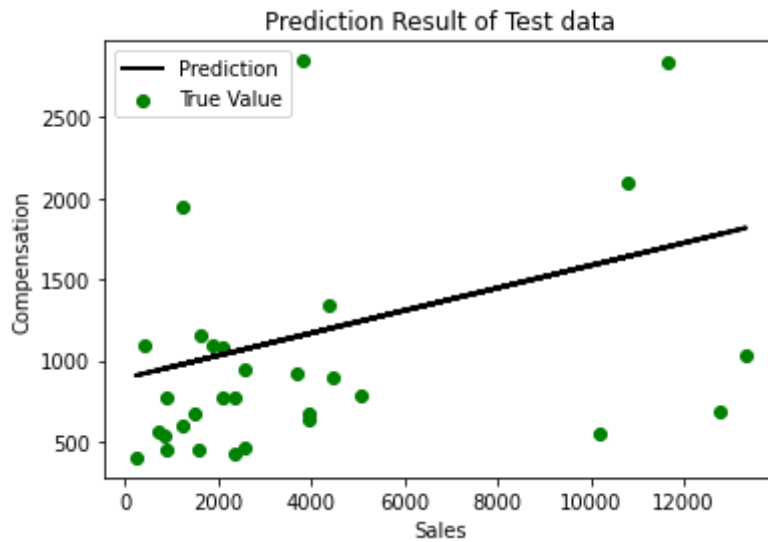
```
plt.show()
```



```
#Prediction Result of Testing data
```

```
Y test pred = model.predict(X test)
```

```
plt.figure()
plt.scatter(X_test, Y_test, color='green', label='True Value')
plt.plot(X_test, Y_test_pred, color='black', linewidth=2, label='Prediction')
plt.xlabel("Sales")
plt.ylabel("Compensation")
plt.title('Prediction Result of Test data')
plt.legend()
plt.show()
```



```
print("Mean squared error =", round(sm.mean_squared_error(Y_test, Y_test_pred), 2))
print("Explain variance score =", round(sm.explained_variance_score(Y_test, Y_test_pred), 2))
print("R2 score =", round(sm.r2_score(Y_test, Y_test_pred), 2))
```

```
Mean squared error = 382854.66
Explain variance score = 0.12
R2 score = 0.05
```