```python
import numpy as np
import pandas as pd
```

```python
import os
for dirname, _, filenames in os.walk('/kaggle/input'):
    for filename in filenames:
        print(os.path.join(dirname, filename))
```

```
/kaggle/input/automobile-dataset/Automobile_data.csv
```

```python
df=pd.read_csv('/kaggle/input/automobile-dataset/Automobile_data.csv')
```

```python
df
```

```python
df.head(5)
```

```python
df.info()
```

```python
list_cate=[]
list_num=[]
for i in list(df.columns):

    if df[i].dtype=='object':

        list_cate.append(i)
    else:
        list_num.append(i)
```

```python
list_cate
```

```python
df['normalized-losses'].unique()
```

```python
df['normalized-losses']=df['normalized-losses'].str.replace('?','145')
```

```python
df['normalized-losses']=df['normalized-losses'].astype('int64')
```

```python
df['make'].unique()
```

```python
df['fuel-type'].unique()
```

```python
df['bore'].unique()
```

```python
df['bore']=df['bore'].str.replace('?','3.35')
```

```python
df['bore']=df['bore'].astype('float')
```

```python
df['stroke'].unique()
```

```python
df['stroke']=df['stroke'].str.replace('?','3.15')
```

```python
df['stroke']=df['stroke'].astype('float')
```

```python
df['horsepower'].unique()
```

```python
df['horsepower']=df['horsepower'].str.replace('?','130')
```

```python
df['horsepower']=df['horsepower'].astype('int64')
```

```python
df['peak-rpm'].unique()
```

```python
df['peak-rpm']=df['peak-rpm'].str.replace('?','4500')
df['peak-rpm']=df['peak-rpm'].astype('int64')
```

```python
df['price']=df['price'].str.replace('?','16000')
df['price']=df['price'].astype('int64')
```

```python
df.info()
```

```python
#now preprocess the categorical features
from sklearn.preprocessing import LabelEncoder
le=LabelEncoder()
for i in list(df.columns):
    if df[i].dtype=='object':
        df[i]=le.fit_transform(df[i])
```

```python
df
```

```python
y=df['price']
x=df.drop(['price'],axis=1)
```

```python
from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test=train_test_split(x,y,random_state=0,test_size=0.2)
```

```python
from sklearn.linear_model import LinearRegression
from sklearn.metrics import r2_score
lr=LinearRegression()
lr.fit(x_train,y_train)
pred_1=lr.predict(x_test)
score_1=r2_score(y_test,pred_1)
```

```python
score_1
```

```python
import matplotlib.pyplot as plt
import seaborn as sns
sns.scatterplot(x=y_test,y=pred_1)
```

```python
#its a decent fit line
```

```python
from sklearn.ensemble import RandomForestRegressor
rfg=RandomForestRegressor()
rfg.fit(x_train,y_train)
pred_2=rfg.predict(x_test)
score_2=r2_score(y_test,pred_2)
```

```python
score_2
```

```python
sns.scatterplot(x=y_test,y=pred_2)
```

```python
#slightly better than previous one
```

```python
from sklearn.ensemble import GradientBoostingRegressor
gbr=GradientBoostingRegressor()
gbr.fit(x_train,y_train)
pred_3=gbr.predict(x_test)
score_3=r2_score(y_test,pred_3)
```

```python
score_3
```

```python
sns.scatterplot(x=y_test,y=pred_3)
```

```python
from sklearn.svm import SVR
svm=SVR()
svm.fit(x_train,y_train)
pred_4=svm.predict(x_test)
score_4=r2_score(y_test,pred_4)
```

```
score_4
```

```
sns.scatterplot(x=y_test,y=pred_4)
```

```
#svm regressor gives the worst fit line
#from all the models randomforest gives the best fit line
```