



Oasis Infobyte Internship, May- 2023

Mayuri Arun Pathak, Data Science Intern

Car Price Prediction With Machine Learning

- Objective: To predict car price with machine learning.
- Dataset: <https://raw.githubusercontent.com/amankharwal/Website-data/master/CarPrice.csv>
- Dataset description:

- Steps we follow:
- 1)Import required libraries
- 2)Data Collection and Processing
- 3)Encoding the Categorical Data:
- 4)Splitting the data and Target
- 5)Splitting Training and Test data
- 6)Model Training
- 7)Model Evaluation
- 8)Visualize the actual prices and Predicted prices

Import required Libraries

```
In [2]: import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.linear_model import Lasso
from sklearn import metrics
```

Data Collection and Processing

```
In [30]: #loading the dataset
car_dataset = pd.read_csv('https://raw.githubusercontent.com/amankharwal/Website-data/master/CarPrice.csv')
```

```
In [31]: # inspecting the first 5 rows of the dataframe
car_dataset.head()
```

aspiration	0
doornumber	0
carbody	0
drivewheel	0
enginelocation	0
wheelbase	0
carlength	0
carwidth	0
carheight	0
curbweight	0
enginestype	0
cylindernumber	0
enginesize	0
fuelsystem	0
boreratio	0
stroke	0
compressionratio	0
horsepower	0

```
In [32]: # checking the number of rows and columns
car_dataset.shape
```

(205, 26)

```
In [33]: # getting some information about the dataset
car_dataset.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 205 entries, 0 to 204
Data columns (total 26 columns):
 # Column Non-Null Count Dtype
---  ---
0 car_id 205 non-null int64
1 symboling 205 non-null int64
2 CarName 205 non-null object
3 fueltype 205 non-null object
4 aspiration 205 non-null object
5 doornumber 205 non-null object
6 carbody 205 non-null object
7 drivewheel 205 non-null object
8 enginelocation 205 non-null object
9 wheelbase 205 non-null float64
10 carlength 205 non-null float64
11 carwidth 205 non-null float64
12 carheight 205 non-null float64
13 curbweight 205 non-null int64
14 enginetype 205 non-null object
15 cylindernumber 205 non-null object
16 enginesize 205 non-null int64
17 fuelsystem 205 non-null object
18 boreratio 205 non-null float64
19 stroke 205 non-null float64
20 compressionratio 205 non-null float64
21 horsepower 205 non-null int64
22 peakrpm 205 non-null int64
23 citympg 205 non-null int64
24 highwaympg 205 non-null int64
25 price 205 non-null float64
dtypes: float64(8), int64(8), object(10)
memory usage: 41.8+ KB
```

```
In [34]: # checking the number of missing values
car_dataset.isnull().sum()
```

```
Out [34]: car_id 0
symboling 0
CarName 0
fueltype 0
aspiration 0
doornumber 0
carbody 0
drivewheel 0
enginelocation 0
wheelbase 0
carlength 0
carwidth 0
carheight 0
curbweight 0
enginetype 0
cylindernumber 0
enginesize 0
fuelsystem 0
boreratio 0
stroke 0
compressionratio 0
horsepower 0
peakrpm 0
citympg 0
highwaympg 0
price 0
dtype: int64
```

```
In [35]: # checking the distribution of categorical data
print(car_dataset.fueltype.value_counts())
print(car_dataset.aspiration.value_counts())
print(car_dataset.doornumber.value_counts())
print(car_dataset.carbody.value_counts())
print(car_dataset.drivewheel.value_counts())
print(car_dataset.enginelocation.value_counts())
print(car_dataset.fuelsystem.value_counts())
print(car_dataset.cylindernumber.value_counts())
print(car_dataset.enginetype.value_counts())
```

```
gas 185
diesel 20
Name: fueltype, dtype: int64
std 168
turbo 37
Name: aspiration, dtype: int64
four 135
two 98
Name: doornumber, dtype: int64
sedan 96
hatchback 70
wagon 25
hardtop 8
convertible 6
Name: carbody, dtype: int64
fwd 120
rwd 76
4wd 9
Name: drivewheel, dtype: int64
front 202
rear 3
Name: enginelocation, dtype: int64
mpfi 94
2bbl 66
l 20
1bbl 11
spdi 9
4bbl 3
rfi 1
spfi 1
Name: fuelsystem, dtype: int64
four 159
six 24
five 11
eight 5
two 4
three 1
twelve 1
Name: cylindernumber, dtype: int64
ohc 148
ohcf 15
ohcv 13
dohc 12
l 12
rotor 4
dohcv 1
Name: enginetype, dtype: int64
```

Encoding the Categorical Data:

```
In [36]: #encoding "fueltype" column
car_dataset.replace({'fueltype':{'gas':0,'diesel':1}},inplace=True)

#encoding "aspiration" column
car_dataset.replace({'aspiration':{'std':0,'turbo':1}},inplace=True)

#encoding "doornumber" column
car_dataset.replace({'doornumber':{'four':1,'two':1}},inplace=True)

#encoding "carbody" column
car_dataset.replace({'carbody':{'sedan':0,'hatchback':1,'wagon':2,'hardtop':3,'convertible':4}},inplace=True)

#encoding "drivewheel" column
car_dataset.replace({'drivewheel':{'fwd':0,'rwd':1,'4wd':2}},inplace=True)

#encoding "enginelocation" column
car_dataset.replace({'enginelocation':{'front':0,'rear':1}},inplace=True)

#encoding "fuelsystem" column
car_dataset.replace({'fuelsystem':{'mpfi':0,'2bbl':1,'l':2,'1bbl':3,'spdi':3,'4bbl':4,'rfi':5,'spfi':6}},inplace=True)

#encoding "cylindernumber" column
car_dataset.replace({'cylindernumber':{'four':0,'six':1,'five':2,'eight':3,'two':4,'three':5,'twelve':6}},inplace=True)

#encoding "enginetype" column
car_dataset.replace({'enginetype':{'ohc':0,'ohcf':1,'ohcv':2,'dohc':3,'l':4,'rotor':5,'dohcv':6}},inplace=True)
```

```
car_dataset.head()
```

202	0	3.58	2.87	8.8	134	5500
203	2	3.01	3.40	23.0	106	4800
204	0	3.78	3.15	9.5	114	5400
citympg	highwaympg					
0	21	27				
1	21	27				
2	19	26				
3	24	30				
4	18	22				
...				
200	23	28				
201	19	25				
202	18	23				
203	26	27				
204	19	25				

Splitting the data and Target

```
In [38]: X = car_dataset.drop(['CarName','price'],axis=1)
Y = car_dataset['price']
```

```
In [39]: print(X)
```

Model Training

1. Linear Regression

```
# loading the linear regression model
lin_reg_model = LinearRegression()

lin_reg_model.fit(X_train,Y_train)

LinearRegression()
```

Model Evaluation

```
# prediction on Training data
training_data_prediction = lin_reg_model.predict(X_train)

# R squared Error
error_score = metrics.r2_score(Y_train, training_data_prediction)
print("R squared Error : ", error_score)

R squared Error : 0.9027472016097514
```

Visualize the actual prices and Predicted prices

```
plt.xlabel("Actual Price")
plt.ylabel("Predicted Price")
plt.title("Actual Prices vs Predicted Prices")
plt.show()
```

```
# R squared Error
error_score = metrics.r2_score(y_test, test_data_prediction)
print("R squared Error = ", error_score)
```

Splitting Training and Test data

```
In [41]: X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size = 0.1, random_state=2)
```

Model Training

1. Linear Regression

```
In [42]: # loading the linear regression model
lin_reg_model = LinearRegression()
```

```
In [43]: lin_reg_model.fit(X_train,Y_train)
```

```
Out [43]: LinearRegression()
```

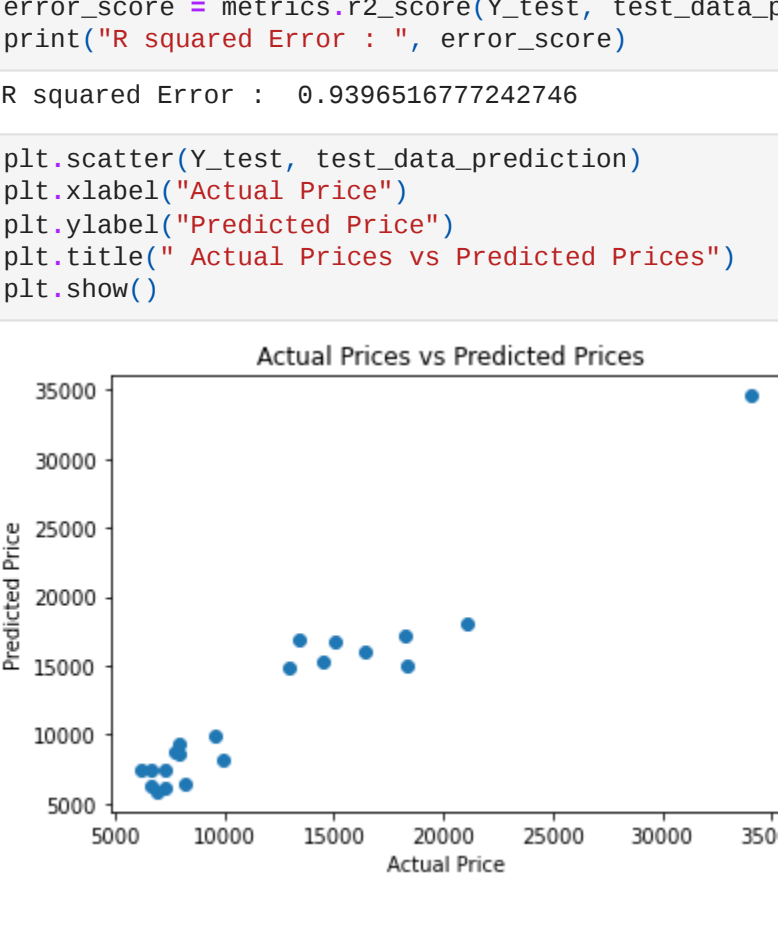
Model Evaluation

```
In [44]: # prediction on Training data
training_data_prediction = lin_reg_model.predict(X_train)
```

```
In [45]: # R Squared Error
error_score = metrics.r2_score(Y_train, training_data_prediction)
print("R squared Error : ", error_score)
R squared Error : 0.9027472016097514
```

Visualize the actual prices and Predicted prices

```
In [46]: plt.scatter(Y_train, training_data_prediction)
plt.xlabel("Actual Price")
plt.ylabel("Predicted Price")
plt.title("Actual Prices vs Predicted Prices")
plt.show()
```



```
In [47]: # prediction on Training data
test_data_prediction = lin_reg_model.predict(X_test)
```

```
In [48]: # R squared Error
error_score = metrics.r2_score(Y_test, test_data_prediction)
print("R squared Error : ", error_score)
R squared Error : 0.9395462837947815
```

```
In [49]: plt.scatter(Y_test, test_data_prediction)
plt.xlabel("Actual Price")
plt.ylabel("Predicted Price")
plt.title("Predicted Price")
plt.title("Actual Prices vs Predicted Prices")
plt.show()
```


- 1. Lasso Regression

```
In [50]: # loading the linear regression model
lass_reg_model = Lasso()
```

```
In [51]: lass_reg_model.fit(X_train,Y_train)
```

```
Out [51]: Lasso()
```

Model Evaluation

```
In [52]: # prediction on Training data
training_data_prediction = lass_reg_model.predict(X_train)
```

```
In [53]: # R squared Error
error_score = metrics.r2_score(Y_train, training_data_prediction)
print("R squared Error : ", error_score)
R squared Error : 0.9027278059187078
```

Visualize the actual prices and Predicted prices

```
In [54]: plt.scatter(Y_train, training_data_prediction)
plt.xlabel("Actual Price")
plt.ylabel("Predicted Price")
plt.title("Predicted Price")
plt.title("Actual Prices vs Predicted Prices")
plt.show()
```



```
In [55]: # prediction on Training data
test_data_prediction = lass_reg_model.predict(X_test)
```

```
In [56]: # R squared Error
error_score = metrics.r2_score(Y_test, test_data_prediction)
print("R squared Error : ", error_score)
R squared Error : 0.9396536777242746
```

```
In [57]: plt.scatter(Y_test, test_data_prediction)
plt.xlabel("Actual Price")
plt.ylabel("Predicted Price")
plt.title("Predicted Price")
plt.title("Actual Prices vs Predicted Prices")
plt.show()
```


Suggestions are always welcome!!

Thank You!