

CAR PRICE PREDICTION

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn import datasets, linear_model
```

```

import seaborn as sns
from sklearn import datasets, linear_model
from sklearn.preprocessing import LabelEncoder, OneHotEncoder, MinMaxScaler
from scipy.stats import shapiro
from scipy.stats import anderson
from scipy.stats import normaltest
from scipy.stats import norm
from scipy.stats import skew
import warnings
warnings.filterwarnings('ignore')

In [2]: df = pd.read_csv("C:/Users/Naailcha CV/OneDrive/Documents/CarPrice_Assignment.csv")
df

Out[2]:
   car_id  symboling  CarName  fueltype  aspiration  doornumber  carbody  drivewheel  enginelocation  wheelbase  _engineize  fuelsystem  boreratio  stroke  compressionratio  horsepower  peakrpm  citympg  highwaympg
0      0         1      3  alfa-romeo  guila      gas      std      two  convertible      rwd      front      88.6  ...      130      mpg      3.47  2.68      9.0      111  5000      21
1      1         2      3  alfa-romeo  romeo      gas      std      two  convertible      rwd      front      88.6  ...      130      mpg      3.47  2.68      9.0      111  5000      21
2      3         3      1  alfa-romeo  romeo      gas      std      two  hatchback      rwd      front      94.5  ...      152      mpg      2.68  3.47      9.0      154  5000      19
3      4         2      2  audi-100  ls      gas      std      four  sedan      fwd      front      99.8  ...      109      mpg      3.19  3.40      10.0      102  5500      24
4      5         2      2  audi-100  ls      gas      std      four  sedan      fwd      front      99.4  ...      136      mpg      3.19  3.40      8.0      115  5500      18
...
200  201        -1      1  volvo  145e  gas      std      four  sedan      rwd      front      109.1  ...      141      mpg      3.78  3.15      9.5      114  5400      23
201  202        -1      1  volvo  144e  gas      turbo      four  sedan      rwd      front      109.1  ...      141      mpg      3.78  3.15      8.7      160  5300      19
202  203        -1      1  volvo  244d  gas      std      four  sedan      rwd      front      109.1  ...      173      mpg      3.58  2.87      8.8      134  5500      18
203  204        -1      1  volvo  246  diesel  turbo      four  sedan      rwd      front      109.1  ...      145      id3      3.01  3.40      23.0      106  4800      26
204  205        -1      1  volvo  264d  gas      turbo      four  sedan      rwd      front      109.1  ...      141      mpg      3.78  3.15      9.5      114  5400      19

205 rows x 26 columns

ATTRIBUTES INFO

car_id: Unique ID # for each car

Symboling: Actuarial assessment of risk of the car (numerical values where -3 is safe, +3 is risky)

CarName: Make and model of car

Fueltype: Car Fuel Type (gas or diesel)

Aspiration: Car Aspiration (standard or turbo)

Doornumber: Number of car doors (two or four)

Carbody: Type of car body (convertible, sedan, hatchback, wagon, or hardtop)

Drivewheel: Car drive wheel (rear wheel drive, 4 wheel drive or front wheel drive)

EngineLocation: Location of car engine (front or rear)

Wheelbase: Car wheel base in inches

CarLength: Car length in inches

CarWidth: Car width in inches

CarHeight: Car height in inches

Curbweight: Car weight in pounds

Enginetype: Car engine type (dohc, dohcv, l, ohc, ohcv, ohcv, or rotor)

Cylindernumber: Number of car cylinders (two, three, four, five, six, eight, or twelve)

EngineSize: Size of engine (numerical values of cubic inches)

Fuelsystem: Type of car fuel system (1bb1, 2 bb1, 4 bb1, id1, mfi, mfg1, spdi, or spfi)

BoreRatio: Car Bore-Stroke Ratio is the ratio between the dimensions of the engine cylinder bore diameter to its piston stroke-length

Stroke: Car strokes (numerical value in strokes)

Compressionratio: Car compression ratio (ratio between the volume of the cylinder with the piston in the bottom position, Vbottom (largest volume), and in the top position, Vtop (smallest volume))

Horsepower: Car horsepower (numerical values of horsepower)

Peakrpm: Car peak RPM (revolutions per minute)

Citympg: Car city MPG (miles per gallon)

Highwaympg: Car highway MPG (miles per gallon)

Price: total price of car in dollars

In [3]: df.head()

Out[3]:
   car_id  symboling  CarName  fueltype  aspiration  doornumber  carbody  drivewheel  enginelocation  wheelbase  _engineize  fuelsystem  boreratio  stroke  compressionratio  horsepower  peakrpm  citympg  highwaympg
0      0         1      3  alfa-romeo  guila      gas      std      two  convertible      rwd      front      88.6  ...      130      mpg      3.47  2.68      9.0      111  5000      21
1      1         2      3  alfa-romeo  romeo      gas      std      two  convertible      rwd      front      88.6  ...      130      mpg      3.47  2.68      9.0      111  5000      21
2      3         3      1  alfa-romeo  romeo      gas      std      two  hatchback      rwd      front      94.5  ...      152      mpg      2.68  3.47      9.0      154  5000      19
3      4         2      2  audi-100  ls      gas      std      four  sedan      fwd      front      99.8  ...      109      mpg      3.19  3.40      10.0      102  5500      24
4      5         2      2  audi-100  ls      gas      std      four  sedan      fwd      front      99.4  ...      136      mpg      3.19  3.40      8.0      115  5500      18

5 rows x 26 columns

In [4]: df.tail()

Out[4]:
   car_id  symboling  CarName  fueltype  aspiration  doornumber  carbody  drivewheel  enginelocation  wheelbase  _engineize  fuelsystem  boreratio  stroke  compressionratio  horsepower  peakrpm  citympg  highwaympg
200  201        -1      1  volvo  145e  gas      std      four  sedan      rwd      front      109.1  ...      141      mpg      3.78  3.15      9.5      114  5400      23
201  202        -1      1  volvo  144e  gas      turbo      four  sedan      rwd      front      109.1  ...      141      mpg      3.78  3.15      8.7      160  5300      19
202  203        -1      1  volvo  244d  gas      std      four  sedan      rwd      front      109.1  ...      173      mpg      3.58  2.87      8.8      134  5500      18
203  204        -1      1  volvo  246  diesel  turbo      four  sedan      rwd      front      109.1  ...      145      id3      3.01  3.40      23.0      106  4800      26
204  205        -1      1  volvo  264d  gas      turbo      four  sedan      rwd      front      109.1  ...      141      mpg      3.78  3.15      9.5      114  5400      19

5 rows x 26 columns

In [5]: # Lets drop the column car_id
df.drop("car_id", axis=1, inplace=True)

In [6]: df.shape
Out[6]:
(295, 25)

```

symboling	0
CarName	0
fueltype	0
aspiration	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

this dataset do not have any null values

In [8]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 285 entries, 0 to 284
Data columns (total 29 columns):
 #   Column                Non-Null Count  Dtype
---  -
 0   symboling              285 non-null    int64
 1   CarName                285 non-null    object
 2   fuelType               285 non-null    object
 3   aspiration              285 non-null    object
 4   doorNumber             285 non-null    object
 5   carbody                285 non-null    object
 6   driveWheel             285 non-null    object
 7   engineLocation         285 non-null    object
 8   wheelbase              285 non-null    float64
 9   carLength              285 non-null    float64
10   carWidth               285 non-null    float64
11   carHeight              285 non-null    float64
12   curbWeight             285 non-null    int64
13   engineType             285 non-null    object
14   cylinderNumber         285 non-null    object
15   engineSize             285 non-null    int64
16   fuelSystem             285 non-null    object
17   boreRatio              285 non-null    float64
18   stroke                 285 non-null    float64
19   compressionRatio       285 non-null    float64
20   horsepower             285 non-null    int64
21   peakRpm                285 non-null    int64
22   cityMpg                285 non-null    int64
23   highwayMpg            285 non-null    int64
24   price                  285 non-null    float64
>>>
```

<pre> In [9]: df.describe() Out[9]: </pre>														
count	250000	205000	205000	205000	205000	205000	205000	205000	205000	205000	205000	205000	205000	205000
mean	10.84316	98.76565	73.40626	65.5685	53.74483	255.56584	12.94287	3.39476	3.25415	10.42537	10.141703	51.251161	25.219512	30.71220
std	1.245307	0.621776	12.37739	2.14504	2.44352	52.68024	0.492631	0.27084	0.13597	3.97740	39.544157	47.985654	5.542142	6.88624
min	2.000000	85.650000	14.150000	30.00000	47.850000	148.000000	61.000000	2.500000	2.070000	7.000000	48.000000	41.000000	13.000000	16.000000
25%	0.000000	94.500000	166.300000	60.00000	52.000000	214.500000	97.000000	3.150000	3.110000	8.000000	70.000000	480.000000	19.000000	25.000000
50%	1.000000	97.000000	173.200000	65.00000	54.100000	241.400000	121.000000	3.330000	3.290000	9.000000	95.000000	520.000000	24.000000	30.000000
75%	1.000000	102.400000	183.100000	65.00000	55.500000	235.500000	141.000000	3.390000	3.410000	9.000000	116.000000	550.000000	30.000000	34.000000
max	3.000000	120.000000	218.000000	72.300000	59.800000	406.000000	328.000000	3.940000	4.170000	13.200000	288.000000	660.000000	49.000000	54.000000
<pre> In [10]: checking for duplicates Out[10]: </pre>														
<pre> In [11]: Out[11]: </pre>														

```
0 rows x 25 columns

In [11]: #Exploring the continuous features of the dataset
num_col = df.select_types(exclude = 'object').columns
num_col

In [11]: Index(['symboling', 'wheelbase', 'carlength', 'carwidth', 'carheight',
              'curbweight', 'engsize', 'boreratio', 'stroke', 'compressionratio',
              'horsepower', 'peakrpm', 'citympg', 'highwaympg', 'price'],
              dtype='object')

In [12]: #Checking for linear relationship between the features and target
plt.figure(figsize=(8,20))
for i, j in enumerate(num_col):
    plt.subplot(5, 2, i+1)
```

The figure consists of three side-by-side scatter plots, each showing the relationship between a different feature and the price of a house. All three plots share the same y-axis, labeled 'price', with a scale from -10,000 to 30,000. Each plot includes a blue regression line and a light blue shaded area representing the confidence interval.

- Left Plot:** The x-axis is labeled 'sqft_living' and ranges from 15 to 55. The data points show a clear negative correlation, with the regression line sloping downwards from left to right.
- Middle Plot:** The x-axis is labeled 'highway_sqft' and ranges from 15 to 55. Similar to the first plot, there is a negative correlation, with the regression line sloping downwards.
- Right Plot:** The x-axis is labeled 'price' and ranges from 5000 to 45000. The data points are tightly clustered around a diagonal line, indicating a very strong positive correlation between the feature and the price.

[illegible]

curbweight	-0.22093	0.77636	0.87728	0.39572	0.0000	0.85054	0.64480	0.16780	0.11362	0.79079	-0.26463	-0.75434	0.77945	0.88326
displacement	-0.00019	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	
horsepower	-0.12901	0.74760	0.96454	0.95150	0.17731	0.64440	0.56774	0.10048	0.00597	0.78797	-0.20648	-0.75434	0.77945	0.88326
stroke	-0.00079	0.10659	0.12653	0.09242	-0.05307	0.0000	0.20319	0.05509	0.0000	0.18610	0.08940	-0.06764	-0.24214	-0.04391
compressionratio	-1.78185	0.24786	0.15844	0.11129	0.26214	0.15362	0.02871	0.00597	0.18610	0.18610	-0.20426	-0.38151	0.23470	0.26201
horsepower	-0.00070	0.33204	0.50265	0.60732	-0.10862	0.0000	0.0000	0.57387	0.0000	0.20426	0.18610	-0.20426	0.23470	0.26201
weight	-0.77006	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	
citympg	-0.05623	0.04744	-0.07699	-0.44704	-0.04460	-0.77144	-0.63568	-0.54632	-0.24145	0.24705	0.18046	-0.11354	0.10000	0.91737
highwaympg	-0.06406	-0.54492	-0.70466	-0.77218	-0.10758	-0.74746	-0.67470	-0.57102	-0.04391	0.26201	-0.77054	-0.54492	-0.91737	0.00000
price	-0.09787	0.11765	0.68320	0.75925	0.11836	0.83006	0.74144	0.55313	0.10443	0.06784	0.0000	-0.08267	0.00000	-0.69799

```
# Lets make the correlation matrix visualise
plt.figure(figsize=(10,10))
sns.heatmap(correlation_matrix, annot=True, cmap="coolwarm", fmt=".2f")
plt.show()
```

	symbiosis	wheelbase	carlength	cowlheight	cowlwidth	cowlarea
symbiosis	1.00	0.36	0.23	0.22	0.11	0.13
wheelbase	0.36	1.00	0.87	0.86	0.87	0.89
carlength	0.23	0.87	1.00	0.84	0.84	0.86
cowlheight	0.22	0.86	0.84	1.00	0.92	0.91
cowlwidth	0.11	0.87	0.84	0.92	1.00	0.98
cowlarea	0.13	0.89	0.86	0.91	0.98	1.00

	curvature	englength	lateralize	gross	compressionratio	herispectrum	peakfreq
curvature	-0.29	0.29	0.00	0.00	0.00	0.00	0.00
englength	-0.01	0.00	0.00	0.00	0.00	0.00	0.00
lateralize	-0.13	0.49	0.11	0.06	0.17	0.00	0.00
gross	-0.01	0.16	0.11	0.16	0.06	0.17	0.00
compressionratio	-0.18	0.25	0.16	0.18	0.26	0.00	0.00
herispectrum	-0.07	0.35	0.55	-0.11	0.76	0.01	0.00
peakfreq	-0.27	0.36	0.29	0.22	0.32	0.27	0.24

[illegible]

```

#stroke
#compression
#peakrpm
#no can also remove CarName as it is insignificant.

df.drop(columns=['symboling', 'carheight', 'stroke', 'compression', 'peakrpm', 'CarName'], axis = 1,
        inplace = True)
df.head()

out[0]:

```

	fueltype	displacement	aspiration	carburetor	driveheel	engine	wheelbase	carlength	carwidth	curbwgt	enginecyl	cylindernumber	engineize	fuelsystem	borestroke	horsepower	citympg	highwaympg	price
0	gas	sd	two	convertible	rd	front	88.6	168.8	64.1	2548	eight	four	130	mp4	3.47	111	21	27	13495.0
1	gas	sd	two	convertible	rd	front	88.6	168.8	64.1	2548	eight	four	130	mp4	3.47	111	21	27	16500.0
2	gas	sd	two	hardback	rd	front	84.5	171.2	65.5	2623	eight	four	130	mp4	2.68	104	19	20	16500.0
3	gas	sd	two	hardback	rd	front	88.6	178.6	66.8	2623	eight	four	130	mp4	2.68	104	24	26	13550.0

```

4 gas          sst          four          sedan          fwd          hork          99.4          176.6          66.4          2824          chx          five          136.6          mpi          3.19          115          18          22          17450.0
In [18]: #Exploring the categorical features of the dataset
cat_col = df.select_dtypes(include = ['object']).columns
df[cat_col].nunique()
Out[18]:
fueltype      2
aspiration    2
doornumber    2
carbody       5
drivewheel    3
engine        2
engine2       7
cylindernumber 7
fuelsystem    8

```

```
dtype: int64

plt.figure(figsize=(20,20))
for i, j in enumerate(cat_cols):
    plt.subplot(3, 3, i+1)
    sns.boxplot(x=j, y='price', data = df)
plt.show()
```

Three bar charts showing the number of cars (n) for different engine types. The first chart shows 'gas' (blue) and 'diesel' (orange). The second chart shows 'at' (blue) and 'turbo' (orange). The third chart shows 'two' (blue) and 'four' (orange). Each chart has a y-axis labeled 'n' ranging from 0 to 10000.

Figure 1 consists of two bar charts. The left chart shows the number of reads (n) for four conditions: blue, orange, green, and purple. The y-axis ranges from 0 to 30,000. The right chart shows the number of reads (p) for the same four conditions. The y-axis ranges from 0 to 20,000. Error bars are present on all bars.

Condition	n (reads)	p (reads)
Blue	~22,000	~19,500
Orange	~11,000	~9,000
Green	~14,500	~11,000
Purple	~22,500	~13,000

Figure 1 consists of six bar charts arranged in a 3x2 grid. Each chart shows a performance metric (AP, F1, AUC, Precision, Recall, F0.5) on the y-axis for different vehicle types or engine types on the x-axis. The proposed method is represented by blue bars, and baseline methods are represented by orange, green, and red bars. Error bars indicate standard deviation.

- Top Left (AP):** Performance for vehicle types: convertible, hatchback, sedan, wagon, pickup, SUV. The proposed method (blue) shows the highest performance across all categories.
- Top Right (F1):** Performance for vehicle types: sedan, #1 (sedan), #2 (SUV). The proposed method (blue) shows the highest performance.
- Middle Left (AUC):** Performance for vehicle types: convertible, hatchback, sedan, wagon, pickup, SUV. The proposed method (blue) shows the highest performance.
- Middle Right (Precision):** Performance for engine types: gas, diesel, electric, hybrid. The proposed method (blue) shows the highest performance.
- Bottom Left (Recall):** Performance for engine types: gas, diesel, electric, hybrid. The proposed method (blue) shows the highest performance.
- Bottom Right (F0.5):** Performance for engine types: gas, diesel, electric, hybrid. The proposed method (blue) shows the highest performance.

Figure 1: Three bar charts showing the distribution of 'group' categories for different 'system' values. The first chart shows 'group' for 'system' values 'dbcs', 'dsv', 'dvc', 'dvc', 'dvc', 'dvc', 'dvc'. The second chart shows 'group' for 'system' values 'four', 'six', 'five', 'three', 'seven', 'two', 'eight'. The third chart shows 'group' for 'system' values 'eight', 'three', 'three', 'three', 'three', 'three', 'three'.

```

    ['gas', 'diesel']
    aspiration
    ['std', 'turbo']
    doornumber
    ['two', 'four']
    carbody
    ['convertible', 'hatchback', 'sedan', 'wagon', 'hardtop']
    driveline
    ['rwd', 'fwd', '4wd']
    enginelocation
    ['front', 'rear']
    enginetype
    ['ohc', 'ohcv', 'ohc', 'l', 'rotor', 'ohcf', 'ohcv']
    cylindernumber
    ['four', 'six', 'five', 'three', 'twelve', 'two', 'eight']
    fuelsystem
    ['four', 'six']

```

One hot encoding

One hot encoding is a technique that we use to represent categorical variables as numerical values in a machine learning model.

```
[ 0]: df = pd.get_dummies(columns = ['fueltype', 'aspiration', 'doornumber', 'driveheel', 'enginelocation'],
df.head()
```

```
[ 0]:
```

	carbody	wheelbase	carlength	carwidth	curbwgt	engine_type	cylindernumber	engineize	fuelsystem	borestroke	horsepower	citympg	highwaympg	price	fueltype_gas	aspiration_turbo	doornumber_two	driveheel_left	enginelocation_front
0	convertible	88.6	100.9	64.1	2540	dohc	four	130	mfi	3.47	111	21	27	13495.0	1	0	0	1	
1	convertible	88.6	100.9	64.1	2540	dohc	four	130	mfi	3.47	111	21	27	16050.0	0	0	0	1	
2	sedan	94.5	112.2	65	2920	di	six	160	tdi	2.68	154	19	25	16950.0	0	0	0	0	1

[illegible]

2	2.0	945	1712	65.5	2023	5.0	3.0	152	5.0	2.68	154	19	26	165000	1	0	1
3	3.0	99.8	176.6	66.2	2337	3.0	2.0	109	5.0	3.19	102	24	30	139500	1	0	0
4	3.0	99.4	176.6	66.4	2624	3.0	1.0	136	5.0	3.19	115	18	22	174500	1	0	0

```
[2]: #model building
from sklearn.model_selection import train_test_split

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

scaler = StandardScaler()
x = scaler.fit_transform(x)

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.3, random_state = 201)
```

```
[22]: from sklearn.linear_model import LinearRegression
      from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score

      lr_model = LinearRegression()
      lr_model.fit(x_train, y_train)

In [22]: LinearRegression()

In [23]: lr_model.fit(x_train, y_train)
         y_pred = lr_model.predict(x_test)

In [24]: pip install yellowbrick
```

[illegible]

```
[20]: lr_pred = lr_model.predict(X_test)
      mae = mean_absolute_error(y_test, lr_pred)
      print('Mean Absolute Error = (mae1:.2f)')
      r2 = lr.score(y_test, lr_pred)
      print('R-Squared Score = (r2:.400:.2f) %')

Mean Absolute Error = 1392.23
R-Squared Score = 78.63 %

In [27]: from yellowbrick.regressor import PredictionError
         plt.figure(figsize=(20,8))
         pred_error = PredictionError(lr_model)

         pred_error.fit(X_train, y_train)
         pred_error.score(X_test, y_test)
         pred_error.show()
```

Prediction Error for Linear Regression

Legend: $R^2 = 0.786$
 - best fit (solid line)
 - identity (dashed line)

The plot shows a positive correlation between actual and predicted values, with the best fit line closely following the identity line.

```
ridge regression(l2 regularization)

ridge is used when you have multiple features and wants to select a better feature correlating to y for that you can use ridge. this can help when features are highly correlated

In [29]: from sklearn.linear_model import Ridge

ridge_model = Ridge(alpha = 50)
ridge_model.fit(X_train, y_train)

ridge_pred = ridge_model.predict(x_test)
mse2 = mean_absolute_error(y_test, ridge_pred)
print('Mean Absolute Error = (mse2, 2%)')

r2 = r2_score(y_test, ridge_pred)
print('R-Squared Score = (r2*100: 2%) %')
```

```

Mean Absolute Error = 1092.13
R-Squared Score = 84.76 %

In [29]: plt.figure(figsize=(20,8))
pred_error = PredictionError(y0g_model)

pred_error.fit(X_train, y_train)
pred_error.score(X_test, y_test)
pred_error.show();


```

Prediction Error for Ridge

$R^2 = 0.848$

test fit

deviation



LASSO REGRESSION (L1 REGULARIZATION)

Lasso regression relies upon the linear regression model but additionally performs a so called L1 regularization, which is a process of introducing additional information in order to prevent overfitting. As a consequence, we can fit a model containing all possible predictors and use lasso to perform variable selection by using a technique that regularizes the coefficient estimates

```
[ 30]: from sklearn.linear_model import Lasso
```

```
lasso_model = Lasso(alpha = 50)
lasso_model.fit(X_train, y_train)

lasso_pred = lasso_model.predict(X_test)
mse = mean_absolute_error(y_test, lasso_pred)
print('Mean Absolute Error = (mse: 27)')

r2 = r2_score(y_test, lasso_pred)
print('R-Squared Score = (r2:100: 27) %')

Mean Absolute Error = 1385.62
R-squared Score = 81.67%
```

```
In [51]: plt.figure(figsize=(20,8))
pred_error = PredictionError(lasso_model)

pred_error.fit(X_train, y_train)
```

```
pred_error = score_X_test, y_test)
pred_error.show()
```

Prediction Error for Lasso

$R^2 = 0.817$

test fit
identity

We can see that three models are a decent fit but Ridge Regression has better generalization.