

# CAR PRICE PREDICTION WITH MACHINE LEARNING

The price of a car depends on a lot of factors like the goodwill of the brand of the car, features of the car, horsepower and the mileage it gives and many more. Car price prediction is one of the major research areas in machine learning. So if you want to learn how to train a car price prediction model

## Import required Modules

```
In [2]: import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

%matplotlib inline
```

## Read CSV file

```
In [3]: car = pd.read_csv('CarPrice_Assignment.csv')
car
car.head()
```

| car_ID | symboling | CarName | fueltpe                  | aspiration | doornumber | carbody         | drivewheel | engineLocation | wheelbase | ... | engineSize | fuelsystem | boreRatio | stroke | compressionRatio | horsepower | peakRpm | cityMpg | highwayMpg |
|--------|-----------|---------|--------------------------|------------|------------|-----------------|------------|----------------|-----------|-----|------------|------------|-----------|--------|------------------|------------|---------|---------|------------|
| 0      | 1         | 3       | alfa-romero guilia       | gas        | std        | two convertible | rwd        | front          | 88.6      | ... | 130        | mpfi       | 3.47      | 2.68   |                  | 9.0        | 111     | 5000    | 21         |
| 1      | 2         | 3       | alfa-romero stelvio      | gas        | std        | two convertible | rwd        | front          | 88.6      | ... | 130        | mpfi       | 3.47      | 2.68   |                  | 9.0        | 111     | 5000    | 21         |
| 2      | 3         | 1       | alfa-romero Quadrifoglio | gas        | std        | two hatchback   | rwd        | front          | 94.5      | ... | 152        | mpfi       | 2.68      | 3.47   |                  | 9.0        | 154     | 5000    | 19         |
| 3      | 4         | 2       | audi 100 ls              | gas        | std        | four sedan      | fwd        | front          | 99.8      | ... | 109        | mpfi       | 3.19      | 3.40   |                  | 10.0       | 102     | 5500    | 24         |
| 4      | 5         | 2       | audi 100ls               | gas        | std        | four sedan      | 4wd        | front          | 99.4      | ... | 136        | mpfi       | 3.19      | 3.40   |                  | 8.0        | 115     | 5500    | 18         |
| ...    | ...       | ...     | ...                      | ...        | ...        | ...             | ...        | ...            | ...       | ... | ...        | ...        | ...       | ...    | ...              | ...        | ...     | ...     | ...        |
| 200    | 201       | -1      | volvo 145e (sw)          | gas        | std        | four sedan      | rwd        | front          | 109.1     | ... | 141        | mpfi       | 3.78      | 3.15   |                  | 9.5        | 114     | 5400    | 23         |
| 201    | 202       | -1      | volvo 144e4a             | gas        | turbo      | four sedan      | rwd        | front          | 109.1     | ... | 141        | mpfi       | 3.78      | 3.15   |                  | 8.7        | 160     | 5300    | 19         |
| 202    | 203       | -1      | volvo 244di              | gas        | std        | four sedan      | rwd        | front          | 109.1     | ... | 173        | mpfi       | 3.58      | 2.87   |                  | 8.8        | 134     | 5500    | 18         |
| 203    | 204       | -1      | volvo 246                | diesel     | turbo      | four sedan      | rwd        | front          | 109.1     | ... | 145        | idi        | 3.01      | 3.40   |                  | 23.0       | 106     | 4800    | 26         |
| 204    | 205       | -1      | volvo 264gl              | gas        | turbo      | four sedan      | rwd        | front          | 109.1     | ... | 141        | mpfi       | 3.78      | 3.15   |                  | 9.5        | 114     | 5400    | 19         |

205 rows × 26 columns

```
In [4]: car.head()
```

| car_ID | symboling | CarName | fueltpe                  | aspiration | doornumber | carbody         | drivewheel | engineLocation | wheelbase | ... | engineSize | fuelsystem | boreRatio | stroke | compressionRatio | horsepower | peakRpm | cityMpg | highwayMpg |    |
|--------|-----------|---------|--------------------------|------------|------------|-----------------|------------|----------------|-----------|-----|------------|------------|-----------|--------|------------------|------------|---------|---------|------------|----|
| 0      | 1         | 3       | alfa-romero guilia       | gas        | std        | two convertible | rwd        | front          | 88.6      | ... | 130        | mpfi       | 3.47      | 2.68   |                  | 9.0        | 111     | 5000    | 21         | 27 |
| 1      | 2         | 3       | alfa-romero stelvio      | gas        | std        | two convertible | rwd        | front          | 88.6      | ... | 130        | mpfi       | 3.47      | 2.68   |                  | 9.0        | 111     | 5000    | 21         | 27 |
| 2      | 3         | 1       | alfa-romero Quadrifoglio | gas        | std        | two hatchback   | rwd        | front          | 94.5      | ... | 152        | mpfi       | 2.68      | 3.47   |                  | 9.0        | 154     | 5000    | 19         | 26 |
| 3      | 4         | 2       | audi 100 ls              | gas        | std        | four sedan      | fwd        | front          | 99.8      | ... | 109        | mpfi       | 3.19      | 3.40   |                  | 10.0       | 102     | 5500    | 24         | 30 |
| 4      | 5         | 2       | audi 100ls               | gas        | std        | four sedan      | 4wd        | front          | 99.4      | ... | 136        | mpfi       | 3.19      | 3.40   |                  | 8.0        | 115     | 5500    | 18         | 22 |

5 rows × 26 columns

```
In [5]: car.tail()
```

| car_ID | symboling | CarName | fueltpe         | aspiration | doornumber | carbody    | drivewheel | engineLocation | wheelbase | ... | engineSize | fuelsystem | boreRatio | stroke | compressionRatio | horsepower | peakRpm | cityMpg | highwayMpg |    |
|--------|-----------|---------|-----------------|------------|------------|------------|------------|----------------|-----------|-----|------------|------------|-----------|--------|------------------|------------|---------|---------|------------|----|
| 200    | 201       | -1      | volvo 145e (sw) | gas        | std        | four sedan | rwd        | front          | 109.1     | ... | 141        | mpfi       | 3.78      | 3.15   |                  | 9.5        | 114     | 5400    | 23         | 28 |
| 201    | 202       | -1      | volvo 144e4a    | gas        | turbo      | four sedan | rwd        | front          | 109.1     | ... | 141        | mpfi       | 3.78      | 3.15   |                  | 8.7        | 160     | 5300    | 19         | 25 |
| 202    | 203       | -1      | volvo 244di     | gas        | std        | four sedan | rwd        | front          | 109.1     | ... | 173        | mpfi       | 3.58      | 2.87   |                  | 8.8        | 134     | 5500    | 18         | 23 |
| 203    | 204       | -1      | volvo 246       | diesel     | turbo      | four sedan | rwd        | front          | 109.1     | ... | 145        | idi        | 3.01      | 3.40   |                  | 23.0       | 106     | 4800    | 26         | 27 |
| 204    | 205       | -1      | volvo 264gl     | gas        | turbo      | four sedan | rwd        | front          | 109.1     | ... | 141        | mpfi       | 3.78      | 3.15   |                  | 9.5        | 114     | 5400    | 19         | 25 |

5 rows × 26 columns

```
In [9]: car.shape
```

(285, 26)

```
In [11]: car.size
```

```
Out[11]: 5330
```

## Data Cleaning

```
In [13]: car.drop(columns=["car_ID", "symboling", "carheight", "stroke", "compressionRatio", "peakRpm"], inplace=True)
car.head()
```

| CarName | fueltpe                  | aspiration | doornumber | carbody         | drivewheel | engineLocation | wheelbase | carlength | carwidth | curbweight | engineType | cylindernumber | engineSize | fuelsystem | boreRatio | horsepower | cityMpg | highwayMpg |
|---------|--------------------------|------------|------------|-----------------|------------|----------------|-----------|-----------|----------|------------|------------|----------------|------------|------------|-----------|------------|---------|------------|
| 0       | alfa-romero guilia       | gas        | std        | two convertible | rwd        | front          | 88.6      | 168.8     | 64.1     | 2548       | dohc       | four           | 130        | mpfi       | 3.47      | 111        | 21      |            |
| 1       | alfa-romero stelvio      | gas        | std        | two convertible | rwd        | front          | 88.6      | 168.8     | 64.1     | 2548       | dohc       | four           | 130        | mpfi       | 3.47      | 111        | 21      |            |
| 2       | alfa-romero Quadrifoglio | gas        | std        | two hatchback   | rwd        | front          | 94.5      | 171.2     | 65.5     | 2823       | ohcv       | six            | 152        | mpfi       | 2.68      | 154        | 19      |            |
| 3       | audi 100 ls              | gas        | std        | four sedan      | fwd        | front          | 99.8      | 176.6     | 66.2     | 2337       | ohc        | four           | 109        | mpfi       | 3.19      | 102        | 24      |            |
| 4       | audi 100ls               | gas        | std        | four sedan      | 4wd        | front          | 99.4      | 176.6     | 66.4     | 2824       | ohc        | five           | 136        | mpfi       | 3.19      | 115        | 18      |            |

```
In [15]: car.shape
```

(285, 20)

```
Out[15]: (285, 20)
```

```
In [16]: car.isnull().sum() #glad, no null values
```

```
Out[16]: CarName      0
fueltpe      0
aspiration    0
doornumber    0
carbody       0
drivewheel    0
engineLocation 0
carlength     0
carwidth     0
curbweight    0
engineType    0
cylindernumber 0
engineSize    0
fuelsystem    0
boreRatio     0
horsepower    0
citympg       0
highwaympg    0
price         0
dtype: int64
```

```
In [18]: car.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 285 entries, 0 to 284
Data columns (total 20 columns):
#   Column      Non-Null Count  Dtype
---  ---
0   CarName      285 non-null    object
1   fueltpe      285 non-null    object
2   aspiration    285 non-null    object
3   doornumber    285 non-null    object
4   carbody      285 non-null    object
5   drivewheel   285 non-null    object
6   engineLocation 285 non-null    object
7   wheelbase    285 non-null    float64
8   carlength    285 non-null    float64
9   carwidth     285 non-null    float64
10  curbweight    285 non-null    int64
11  engineType    285 non-null    object
12  cylindernumber 285 non-null    object
13  engineSize    285 non-null    object
14  fuelsystem    285 non-null    object
15  boreRatio     285 non-null    float64
16  horsepower    285 non-null    int64
17  citympg       285 non-null    int64
18  highwaympg    285 non-null    int64
19  price         285 non-null    float64
dtypes: float64(5), int64(5), object(10)
memory usage: 32.2+ KB
```

```
In [17]: car.describe()
```

|       | wheelbase  | carlength  | carwidth   | curbweight  | engineSize | boreRatio  | horsepower | citympg    | highwaympg | price        |
|-------|------------|------------|------------|-------------|------------|------------|------------|------------|------------|--------------|
| count | 285.000000 | 285.000000 | 285.000000 | 285.000000  | 285.000000 | 285.000000 | 285.000000 | 285.000000 | 285.000000 | 285.000000   |
| mean  | 98.756585  | 174.049268 | 65.907805  | 2555.565854 | 126.907317 | 3.329756   | 104.117073 | 25.219512  | 30.751220  | 13276.710571 |
| std   | 6.021776   | 12.337289  | 2.145204   | 520.680204  | 41.642693  | 0.270844   | 39.544167  | 6.542142   | 6.886443   | 7988.852332  |
| min   | 86.600000  | 141.100000 | 60.300000  | 1488.000000 | 61.000000  | 2.540000   | 48.000000  | 13.000000  | 16.000000  | 5118.000000  |
| 25%   | 94.500000  | 166.300000 | 64.100000  | 2145.000000 | 97.000000  | 3.150000   | 70.000000  | 19.000000  | 25.000000  | 7788.000000  |
| 50%   | 97.000000  | 173.200000 | 65.500000  | 2414.000000 | 120.000000 | 3.310000   | 95.000000  | 24.000000  | 30.000000  | 10295.000000 |
| 75%   | 102.400000 | 183.100000 | 66.900000  | 2935.000000 | 141.000000 | 3.580000   | 116.000000 | 30.000000  | 34.000000  | 16503.000000 |
| max   | 120.900000 | 208.100000 | 72.300000  | 4066.000000 | 326.000000 | 3.940000   | 288.000000 | 49.000000  | 54.000000  | 45400.000000 |

```
In [24]: # cat_features = []
for i in [1, 2, 3, 4, 5, 6, 11, 12, 14]:
    print(f'{car.columns[i]} \t {car[car.columns[i]].unique()}')
    cat_features.append(car.columns[i])

cat_features.append("CarName")

fueltpe      ['gas', 'diesel']
aspiration    ['std', 'turbo']
doornumber    ['two', 'four']
carbody       ['convertible', 'hatchback', 'sedan', 'wagon', 'hardtop']
drivewheel    ['rwd', 'fwd']
engineLocation ['front', 'rear']
engineType    ['dohc', 'ohcv', 'ohc', 'i', 'rotor', 'ohcf', 'dohcv']
cylindernumber ['four', 'six', 'five', 'three', 'twelve', 'two', 'eight']
fuelsystem    ['mpfi', '2bbl', 'mfi', '1bbl', 'spfi', '4bbl', 'idi', 'spdi']
```

```
In [20]: car["drivewheel"].value_counts()
```

```
Out[20]: fwd      129
rwd       76
4wd        9
Name: drivewheel, dtype: int64
```

```
In [21]: car["drivewheel"] = car["drivewheel"].replace('4wd', 'fwd')
car["drivewheel"].value_counts()
```

```
Out[21]: fwd      129
rwd       76
Name: drivewheel, dtype: int64
```

## Heatmap

```
In [25]: plt.figure(figsize=(12, 5))
sns.heatmap(car.corr(), cmap="YlOrRd", annot=True, lw=0.5)
plt.show()
```

## Training the Model

Importing required machine learning libraries for training the model

```
In [26]: from sklearn.model_selection import train_test_split
from sklearn import metrics
from sklearn.preprocessing import OneHotEncoder
from sklearn.compose import make_column_transformer
from sklearn.pipeline import make_pipeline
```

```
In [28]: X = car.iloc[:, :19]
y = car.iloc[:, -1]

print(X.shape)
print(y.shape)
```

(285, 19)

(285,)

```
In [30]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.20, random_state=42)
```

## OneHotEncoding

```
In [31]: ohe = OneHotEncoder()
ohe.fit(X[cat_features])

trf1 = make_column_transformer(
    (OneHotEncoder(categories=ohe.categories_), cat_features),
    remainder='passthrough'
)
```

## Linear Regression

```
In [32]: from sklearn.linear_model import LinearRegression

lin_reg = LinearRegression()
pipe_l = make_pipeline(trf1, lin_reg)

pipe_l.fit(X_train, y_train)
```

```
Out[32]: Pipeline(steps=[('columntransformer',
                        ColumnTransformer(remainder='passthrough',
                        transformers=[('onehotencoder',
                                      OneHotEncoder(categories=[array(['diesel', 'gas'], dtype=object),
                                                                array(['std', 'turbo'], dtype=object),
                                                                array(['four', 'two'], dtype=object),
                                                                array(['convertible', 'hardtop', 'hatchback', 'sedan', 'wagon'],
                                                                dtype=object),
                                                                array(['fwd', 'rwd'], dtype=object...
                                      ['volvo 264gl', 'volvo diesel', 'vw dasher', 'vw rabbit'],
                                      dtype=object)]),
                        ['fueltpe', 'aspiration',
                          'doornumber', 'carbody',
                          'drivewheel',
                          'engineLocation',
                          'engineType',
                          'cylindernumber',
                          'fuelsystem', 'CarName',
                          'fueltpe', 'aspiration',
                          'doornumber', 'carbody',
                          'drivewheel',
                          'engineLocation',
                          'engineType',
                          'cylindernumber',
                          'fuelsystem',
                          'CarName']]])),
                  ('linearregression', LinearRegression()))])
```

```
In [34]: y_predLM = pipe_l.predict(X_test)
```

```
In [35]: r2_lm = metrics.r2_score(y_test, y_predLM).round(4)
r2_lm
```

```
Out[35]: 0.619
```

```
In [36]: RMSE_lm = metrics.mean_absolute_error(y_test, y_predLM).round(3)
MAE_lm = np.sqrt(metrics.mean_squared_error(y_test, y_predLM)).round(3)
print(f'MAE = {MAE_lm} \nRMSE = {RMSE_lm}')

MAE = 5484.144
RMSE = 3467.494
```

## Decision Tree Regressor

```
In [37]: from sklearn.tree import DecisionTreeRegressor

dt = DecisionTreeRegressor()
pipe_dt = make_pipeline(trf1, dt)

pipe_dt.fit(X_train, y_train)
```

```
Out[37]: Pipeline(steps=[('columntransformer',
                        ColumnTransformer(remainder='passthrough',
                        transformers=[('onehotencoder',
                                      OneHotEncoder(categories=[array(['diesel', 'gas'], dtype=object),
                                                                array(['std', 'turbo'], dtype=object),
                                                                array(['four', 'two'], dtype=object),
                                                                array(['convertible', 'hardtop', 'hatchback', 'sedan', 'wagon'],
                                                                dtype=object),
                                                                array(['fwd', 'rwd'], dtype=object...
                                      ['volvo 264gl', 'volvo diesel', 'vw dasher', 'vw rabbit'],
                                      dtype=object)]),
                        ['fueltpe', 'aspiration',
                          'doornumber', 'carbody',
                          'drivewheel',
                          'engineLocation',
                          'engineType',
                          'cylindernumber',
                          'fuelsystem', 'CarName',
                          'fueltpe', 'aspiration',
                          'doornumber', 'carbody',
                          'drivewheel',
                          'engineLocation',
                          'engineType',
                          'cylindernumber',
                          'fuelsystem',
                          'CarName']]])),
                  ('decisiontreeregressor', DecisionTreeRegressor())])
```

```
In [38]: y_predDT = pipe_dt.predict(X_test)
```

```
In [39]: r2_dt = metrics.r2_score(y_test, y_predDT).round(4)
r2_dt
```

```
Out[39]: 0.8729
```

```
In [40]: RMSE_dt = metrics.mean_absolute_error(y_test, y_predDT).round(3)
MAE_dt = np.sqrt(metrics.mean_squared_error(y_test, y_predDT)).round(3)
print(f'MAE = {MAE_dt} \nRMSE = {RMSE_dt}')

MAE = 3167.385
RMSE = 1969.346
```

## Random Forest

```
In [41]: from sklearn.ensemble import RandomForestRegressor

rf = RandomForestRegressor(random_state=0)
pipe_rf = make_pipeline(trf1, rf)

pipe_rf.fit(X_train, y_train)
```

```
Out[41]: Pipeline(steps=[('columntransformer',
                        ColumnTransformer(remainder='passthrough',
                        transformers=[('onehotencoder',
                                      OneHotEncoder(categories=[array(['diesel', 'gas'], dtype=object),
                                                                array(['std', 'turbo'], dtype=object),
                                                                array(['four', 'two'], dtype=object),
                                                                array(['convertible', 'hardtop', 'hatchback', 'sedan', 'wagon'],
                                                                dtype=object),
                                                                array(['fwd', 'rwd'], dtype=object...
                                      ['volvo 264gl', 'volvo diesel', 'vw dasher', 'vw rabbit'],
                                      dtype=object)]),
                        ['fueltpe', 'aspiration',
                          'doornumber', 'carbody',
                          'drivewheel',
                          'engineLocation',
                          'engineType',
                          'cylindernumber',
                          'fuelsystem', 'CarName',
                          'fueltpe', 'aspiration',
                          'doornumber', 'carbody',
                          'drivewheel',
                          'engineLocation',
                          'engineType',
                          'cylindernumber',
                          'fuelsystem',
                          'CarName']]])),
                  ('randomforestregressor',
                  RandomForestRegressor(random_state=0))])
```

```
In [42]: y_predRF = pipe_rf.predict(X_test)
```

```
In [43]: r2_rf = metrics.r2_score(y_test, y_predRF).round(4)
r2_rf
```

```
Out[43]: 0.9531
```

```
In [44]: RMSE_rf = metrics.mean_absolute_error(y_test, y_predRF).round(3)
MAE_rf = np.sqrt(metrics.mean_squared_error(y_test, y_predRF)).round(3)
print(f'MAE = {MAE_rf} \nRMSE = {RMSE_rf}')

MAE = 1924.679
RMSE = 1356.174
```

## Conclusion

```
In [46]: pd.DataFrame(data=[['r2_lm', MAE_lm, RMSE_lm], [r2_dt, MAE_dt, RMSE_dt], [r2_rf, MAE_rf, RMSE_rf]],
                      index=['Logistic Regression', 'Decision Tree', 'Random Forest'],
                      columns=['R2 Score', 'MAE', 'RMSE'])
```

```
Out[46]:
```

|                     | R2 Score | MAE      | RMSE     |
|---------------------|----------|----------|----------|
| Logistic Regression | 0.6190   | 5484.144 | 3467.494 |
| Decision Tree       | 0.8729   | 3167.385 | 1969.346 |
| Random Forest       | 0.9531   | 1924.679 | 1356.174 |