

PGPDSE FT Mini-Project – Machine Learning Title: Car Price Prediction

PARTICIPANT LIST:

1. Piyush Khupse
2. Sahil Radke
3. Pratik Haldankar
4. Shubham Rajput
5. Prajakta Ippar
6. Sakshi Sonar

PROBLEM STATEMENT:

A Chinese automobile company Geely Auto aspires to enter the US market by setting up their manufacturing unit there and producing cars locally to give competition to their US and European counterparts. They have contracted an automobile consulting company to understand the factor on which the price of cars depends. Specifically, they want to understand the factors affecting the pricing of cars in the American market, since those may be very different from the Chinese market.

The company wants to know:

- Which variables are significant in predicting the price of a car.
- How well those variables describe the price of a car.

Based on various market surveys, the consulting firm has gathered a large dataset of different types of cars across the American market

ATTRIBUTES:

- Car_ID - Unique ID for each observation.
- Symboling - Its assigned insurance risk rating, value +3 indicates that the auto is risky, 3 that it is pretty safe.
- carCompany - Name of company.
- fueltype - Car fuel type.
- aspiration - Aspiration used in car.
- doornumber - Number of doors in a car.
- carbody - body of car.
- drivewheel - type of drive wheel.

- enginelocation - location of car engine
- wheelbase - Wheelbase of car.
- carlength - length of car.
- carwidth - width of car.
- carheight - height of car.
- curbweight - The weight of a car without occupants or luggage.
- enginetype - type of engine.
- cylindernumber - cylinder placed in the car.
- enginesize - size of car.
- fuelsystem - Fuel system of car.
- boreratio - Boreratio of car.
- stroke - Stroke or volume inside the engine.
- compressionratio - compression ratio of car.
- horsepower - Horsepower
- peakrpm - car peak rpm
- citympg - Mileage in city
- highwaympg - Mileage on highway

OBJECTIVES:

You are required to model the prices of cars with the available independent variables. It will be used by management to understand how exactly the prices vary with the independent variables. They can accordingly manipulate the design of the cars, the business strategy etc. to meet certain price levels. Further, the model will be good for management to understand the pricing dynamics of the new market.

Step 1:

Importing the libraries:

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from sklearn.model_selection import train_test_split
6 from sklearn.linear_model import LinearRegression
7 from sklearn.metrics import mean_squared_error
8 import warnings
9 warnings.filterwarnings('ignore')
```

Step 2:

Reading the Dataset:

```
1 # Importing the dataset
2 car_data = pd.read_csv('CarPrice_Assignment.csv')
3 car_data.head()
```

	car_ID	symboling	CarName	fueltype	aspiration	doornumber	carbody	drivewheel	engine location	wheelbase	...	enginesize	fuelsystem	bore ratio	stroke
0	1	3	alfa-romero giulia	gas	std	two	convertible	rwd	front	88.6	...	130	mpfi	3.47	
1	2	3	alfa-romero stelvio	gas	std	two	convertible	rwd	front	88.6	...	130	mpfi	3.47	
2	3	1	alfa-romero Quadrifoglio	gas	std	two	hatchback	rwd	front	94.5	...	152	mpfi	2.68	
3	4	2	audi 100 ls	gas	std	four	sedan	fwd	front	99.8	...	109	mpfi	3.19	
4	5	2	audi 100ls	gas	std	four	sedan	4wd	front	99.4	...	136	mpfi	3.19	

5 rows × 26 columns

Step 3:

Checking the dimension of the dataset:

```
1 ## Checking the dimensions of the dataset
2 print("Dimensions of the dataset: ", car_data.shape)
```

Dimensions of the dataset: (205, 26)

Checking the datatype of the columns:

```
1 ## Checking the data types of columns
2 print("\nData types of columns:\n", car_data.dtypes)
```

```
Data types of columns:
car_ID          int64
symboling       int64
CarName         object
fueltype        object
aspiration       object
doornumber       object
carbody         object
drivewheel       object
engine location  object
wheelbase       float64
carlength       float64
carwidth        float64
carheight       float64
curbweight      int64
enginetype       object
cylindernumber   object
enginesize      int64
fuelsystem       object
bore ratio      float64
stroke          float64
compressionratio float64
horsepower      int64
peakrpm         int64
citympg         int64
highwaympg      int64
price           float64
dtype: object
```

Checking for missing values:

```
1 ## Checking for missing values
2 print("\nMissing values:\n", car_data.isnull().sum())
```

Missing values:

```
car_ID      0
symboling   0
CarName      0
fueltype     0
aspiration   0
doornumber   0
carbody      0
drivewheel   0
engineloction 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
```

Getting the Descriptive Statistics of the data:

```
1 ## Descriptive statistics
2 print("\nDescriptive statistics:\n", car_data.describe().T)
3
```

Descriptive statistics:

	count	mean	std	min	25%	50%	75%	max
car_ID	205.0	103.000000	59.322565	1.00	52.00	103.00	154.00	205.00
symboling	205.0	0.834146	1.245307	-2.00	0.00	1.00	2.00	3.00
wheelbase	205.0	98.756585	6.021776	86.60	94.50	97.00	102.40	120.90
carlength	205.0	174.049268	12.337289	141.10	166.30	173.20	183.10	208.10
carwidth	205.0	65.907805	2.145204	60.30	64.10	65.50	66.90	72.30
carheight	205.0	53.724878	2.443522	47.80	52.00	54.10	55.50	59.80
curbweight	205.0	2555.565854	520.680204	1488.00	2145.00	2414.00	2935.00	4066.00
enginesize	205.0	126.907317	41.642693	61.00	97.00	120.00	141.00	326.00
boreratio	205.0	3.329756	0.270844	2.54	3.15	3.31	3.58	3.94
stroke	205.0	3.255415	0.313597	2.07	3.11	3.29	3.41	4.17
compressionratio	205.0	10.142537	3.972040	7.00	8.60	9.00	9.40	23.00
horsepower	205.0	104.117073	39.544167	48.00	70.00	95.00	116.00	288.00
peakrpm	205.0	5125.121951	476.985643	4150.00	4800.00	5200.00	5500.00	6600.00
citympg	205.0	25.219512	6.542142	13.00	19.00	24.00	30.00	49.00
highwaympg	205.0	30.751220	6.886443	16.00	25.00	30.00	34.00	54.00
price	205.0	13276.710571	7988.852332	5118.00	7788.00	10295.00	16503.00	45400.00

Step 4:

Creating a new column 'CompanyName' using 'CarName' column:

```
1 # Creating a new column 'CompanyName' using 'CarName' column
2 car_data['CompanyName'] = car_data['CarName'].apply(lambda x: x.split()[0])
3
```

Checking the correctness of data in the 'CompanyName' column:

```
1 # Checking the correctness of data in the 'CompanyName' column
2 car_data.loc[car_data['CompanyName'] == 'vw', 'CompanyName'] = 'volkswagen'
3 car_data.loc[car_data['CompanyName'] == 'vokswagen', 'CompanyName'] = 'volkswagen'
4 car_data.loc[car_data['CompanyName'] == 'porcshe', 'CompanyName'] = 'porsche'
5 car_data.loc[car_data['CompanyName'] == 'toyouta', 'CompanyName'] = 'toyota'
6 car_data.loc[car_data['CompanyName'] == 'Nissan', 'CompanyName'] = 'nissan'
7 car_data.loc[car_data['CompanyName'] == 'maxda', 'CompanyName'] = 'mazda'
```

Checking for duplicate data:

```
1 # Checking for duplicate data
2 print(car_data.duplicated().sum()) # printing the number of duplicate rows
```

0

INFERENCES:

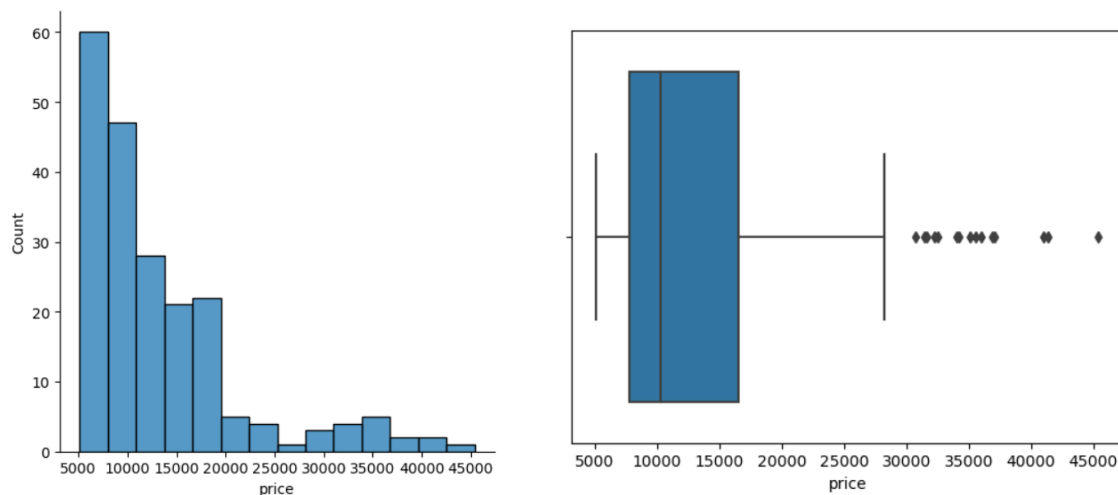
- There are a total of 22 unique car company names in the dataset.
- The 'CompanyName' column looks correct as there are no misspelled or duplicate names.
- There are no duplicate rows in the dataset.

Exploratory Data Analysis

Step 5:

Visualizing the 'price' column using distplot and boxplot:

```
1 # Visualizing the 'price' column using distplot and boxplot
2 sns.distplot(car_data, x='price')
3 plt.show()
4 sns.boxplot(x='price', data=car_data)
5 plt.show()
```



INFERENCE:

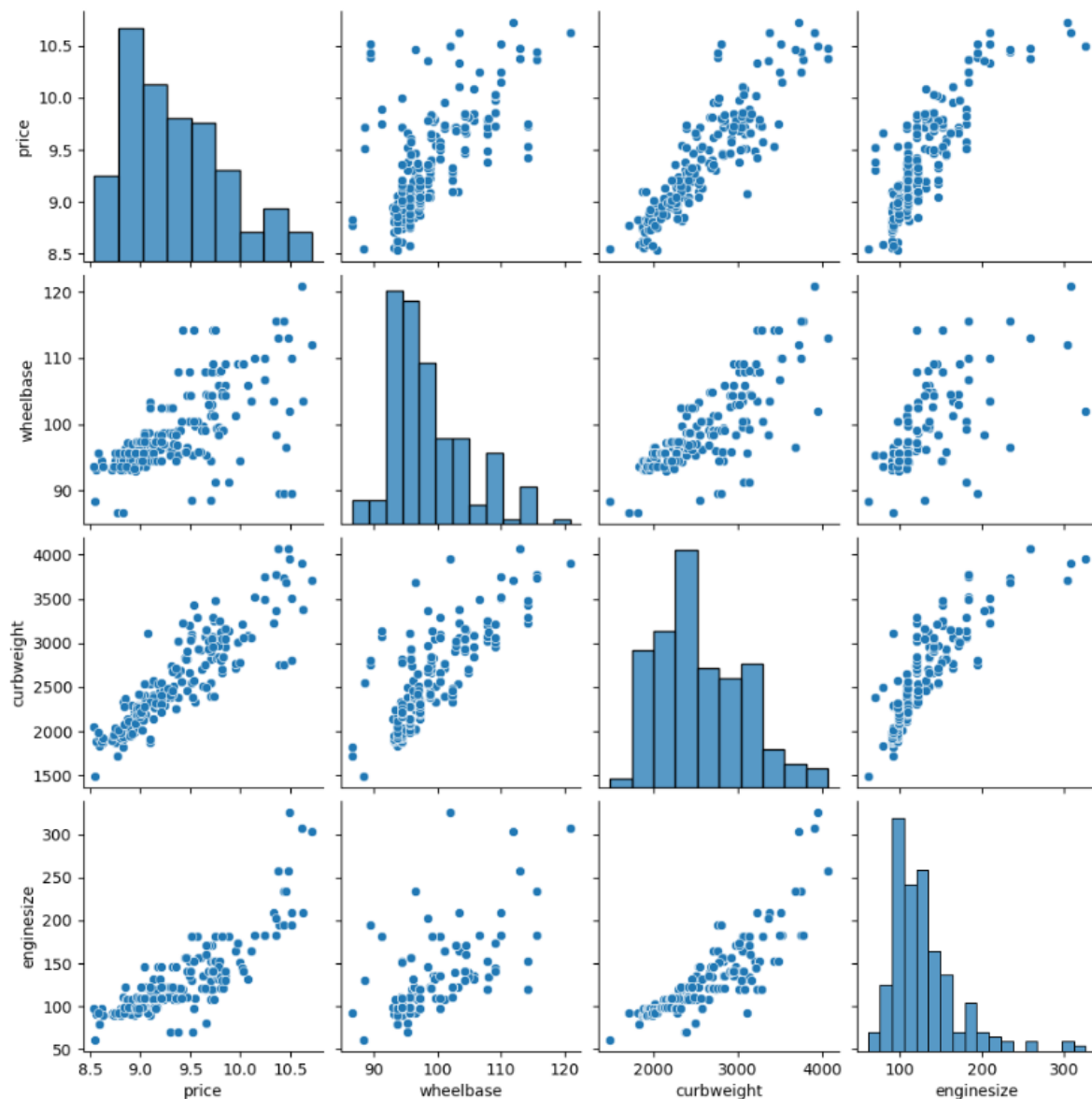
- The 'price' column is right-skewed and has some outliers.

Performing the appropriate transformation to make the target as a gaussian distribution:

```
1 # Performing the appropriate transformation to make the target as a gaussian distribution
2 car_data['price'] = np.log(car_data['price'])
```

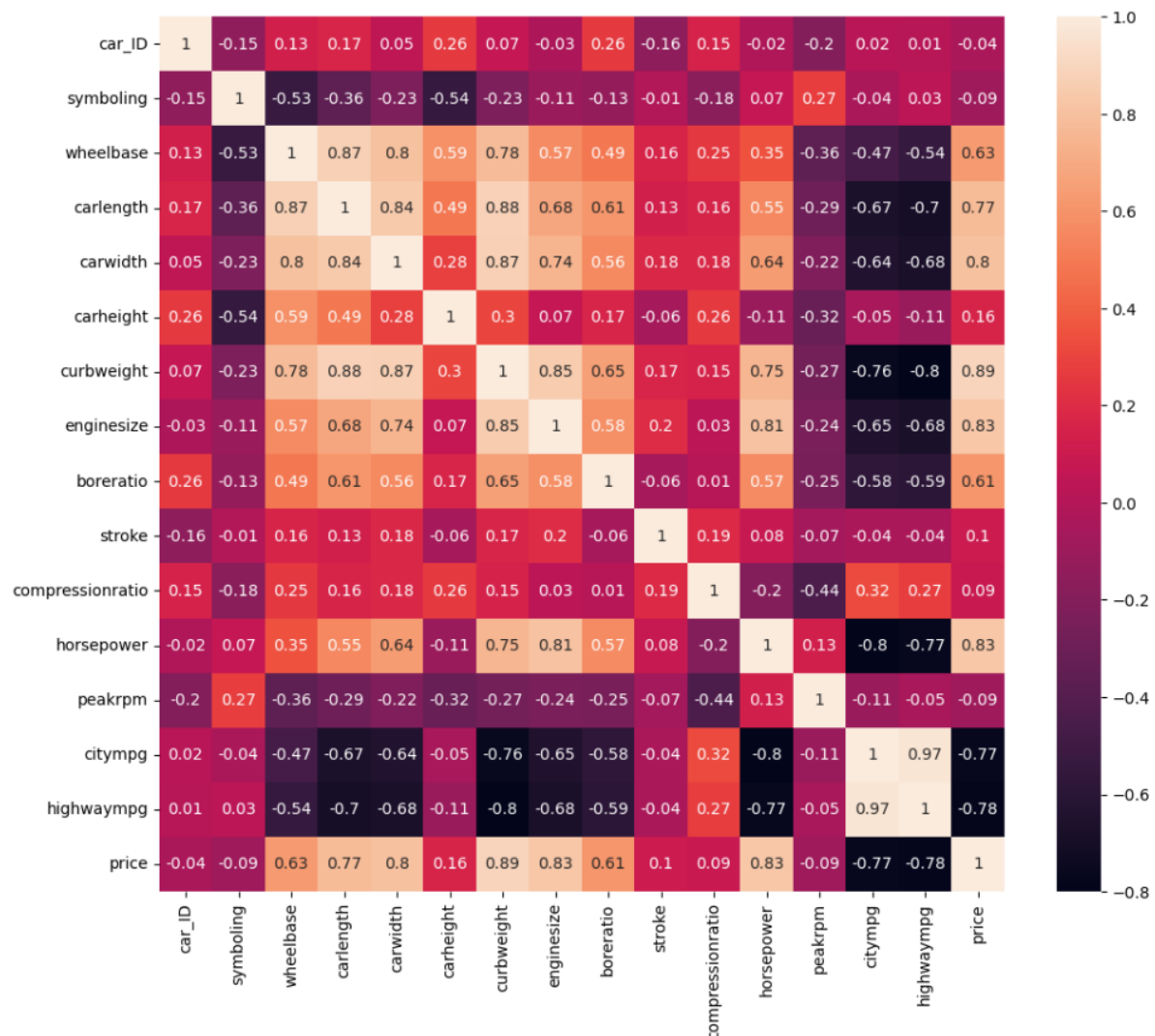
Checking the linear relationship between the dependent variable "Price" and the numerical independent variables:

```
1 # Checking the linear relationship between the dependent variable "Price" and the numerical independent variables
2 sns.pairplot(car_data, vars=['price', 'wheelbase', 'curbweight', 'engineize'])
3 plt.show()
```



Checking the multicollinearity between the correlated independent variables above and Price:

```
1 # Checking the multicollinearity between the correlated independent variables above and Price
2 plt.figure(figsize=(12, 10))
3 correlation_matrix = car_data.corr().round(2)
4 sns.heatmap(data=correlation_matrix, annot=True)
5 plt.show()
```



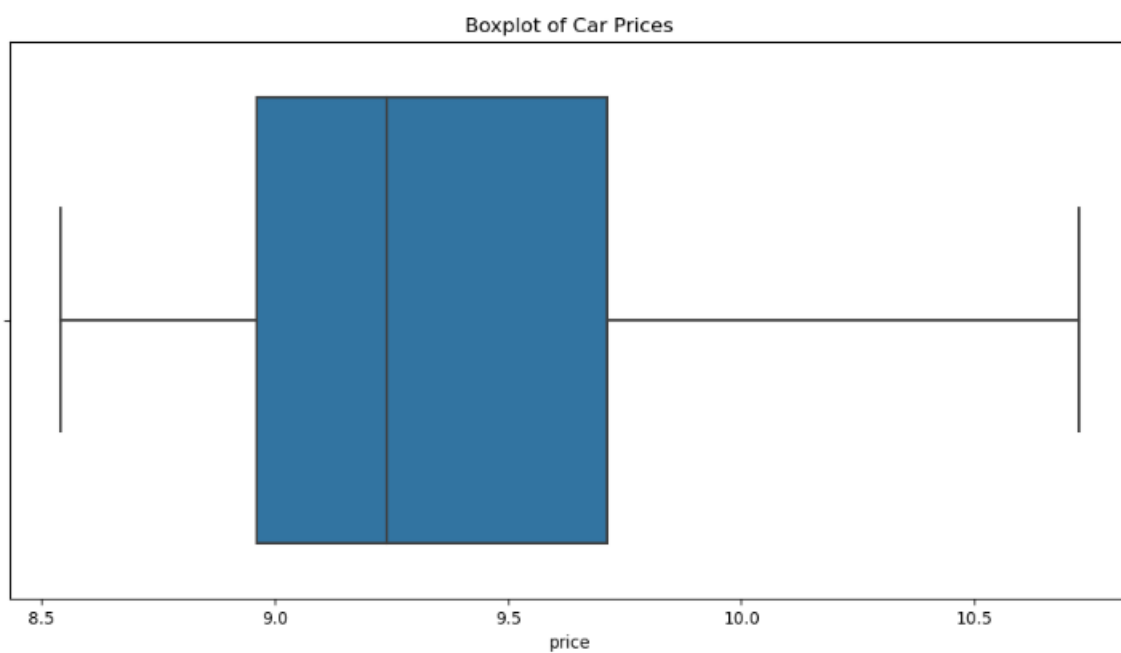
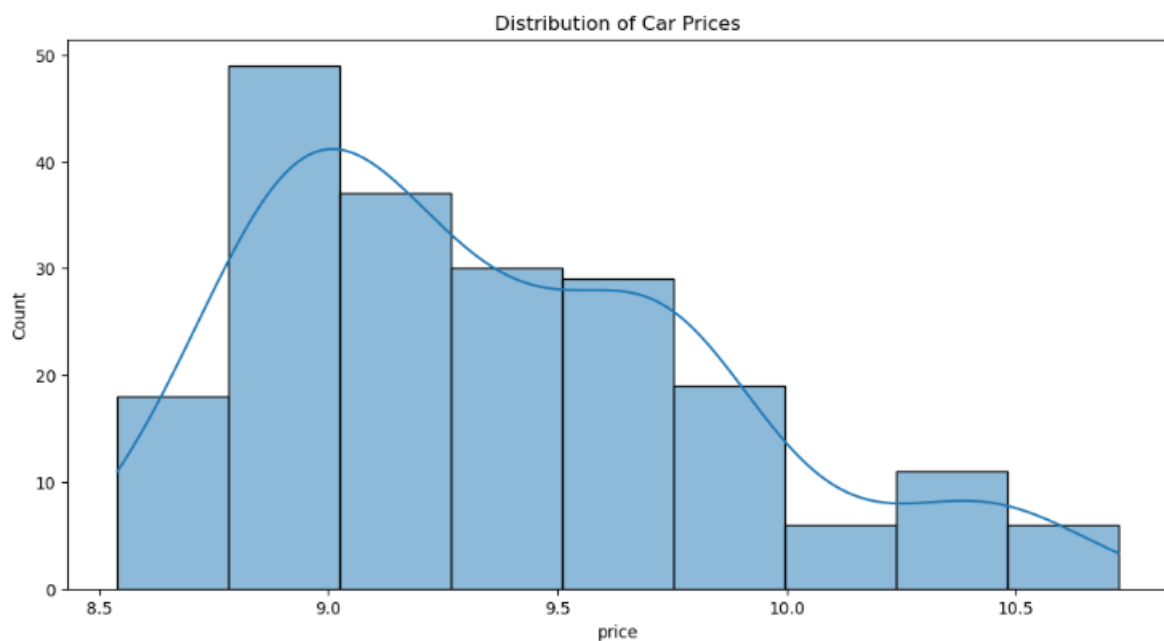
INFERENCES:

- 'wheelbase', 'curbweight', and 'enginesize' have a positive linear relationship with 'price'.
- 'wheelbase' and 'curbweight' are highly correlated with each other.
- 'enginesize' and 'horsepower' are highly correlated with each other.
- 'carlength' and 'carwidth' are highly correlated with 'curbweight'.
- 'citympg' and 'highwaympg' are highly correlated with each other and negatively correlated with 'price'.
- The 'price' column is positively skewed with a long tail, indicating that there are some expensive cars in the dataset.
- The boxplot of car prices shows some outliers, which may need to be removed before training the model.
- The log transformation is applied to the 'price' column to make it a more Gaussian-like distribution.
- The scatterplots show a linear relationship between the price and the numerical independent variables, indicating that these variables can be good predictors of the car price.
- The correlation matrix and pair plot show a high correlation between some of the independent variables, indicating that they may suffer from multicollinearity.
- It is recommended to remove one of the highly correlated variables to avoid this issue.

- Feature engineering is performed by converting the 'cylindernumber' column to a numerical variable.

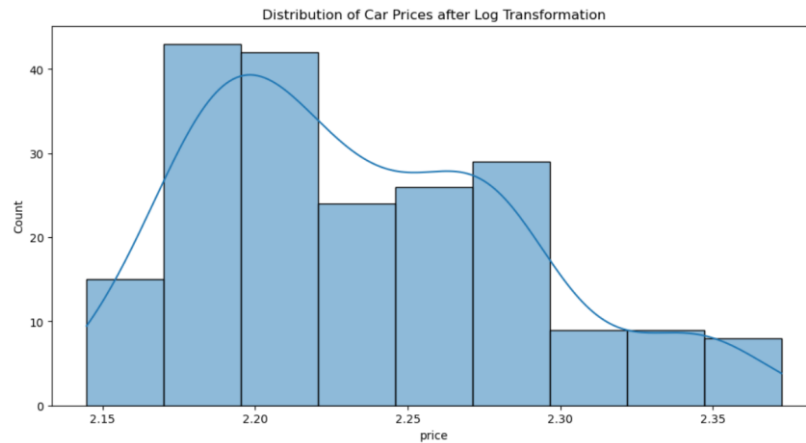
Visualize the distribution of price using a histogram and boxplot:

```
1 # Visualize the distribution of price using a histogram and boxplot
2 plt.figure(figsize=(12,6))
3 sns.histplot(car_data['price'], kde=True)
4 plt.title('Distribution of Car Prices')
5 plt.show()
6
7 plt.figure(figsize=(12,6))
8 sns.boxplot(car_data['price'])
9 plt.title('Boxplot of Car Prices')
10 plt.show()
11
```



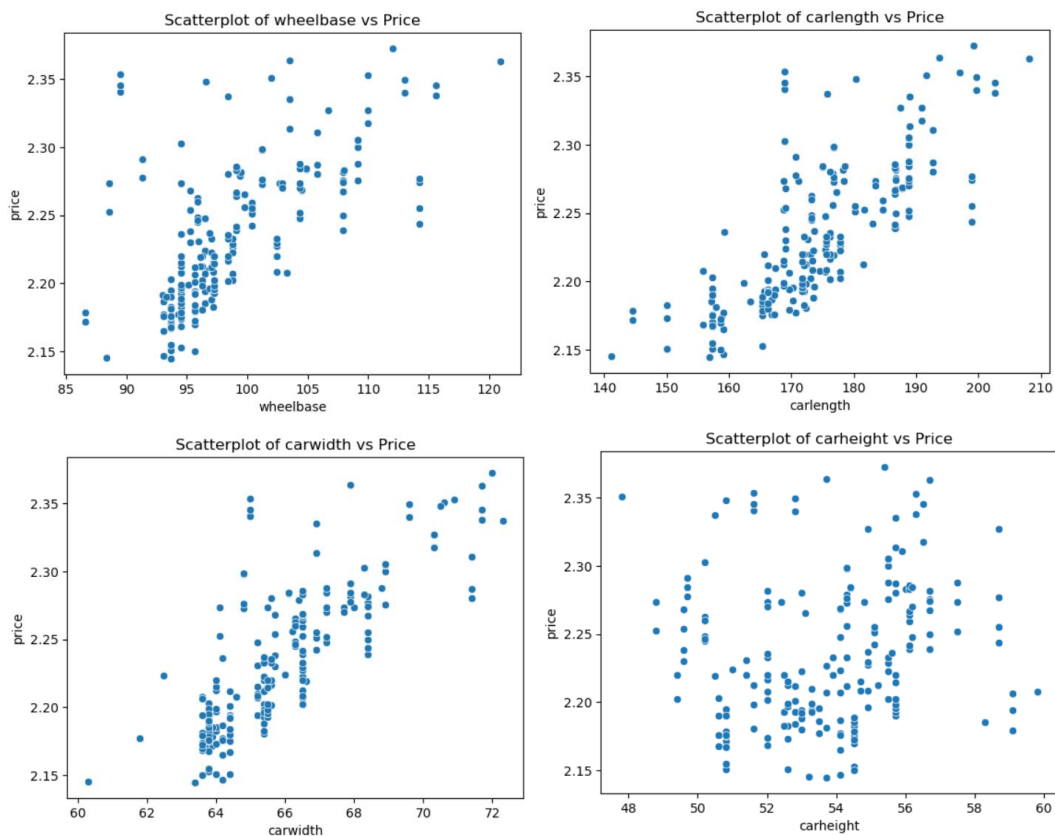
Transforming the target variable using a log transformation to make it normally distributed:

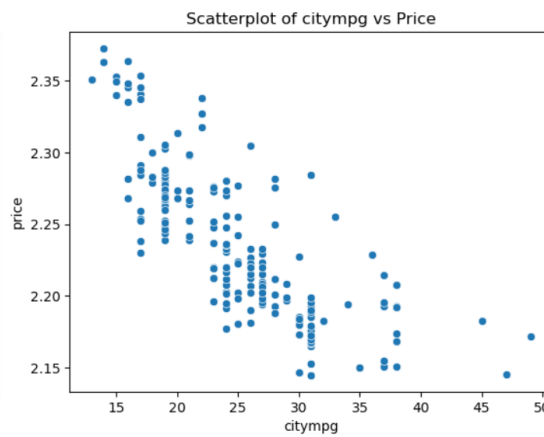
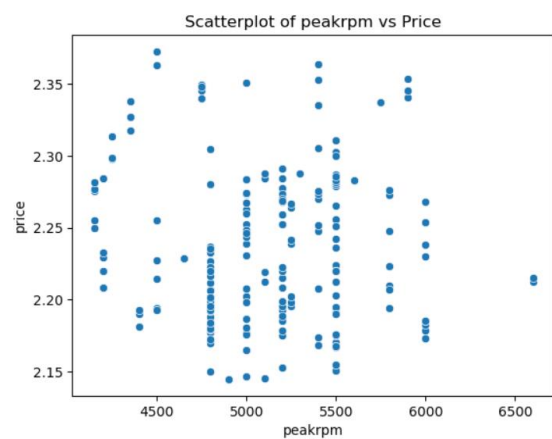
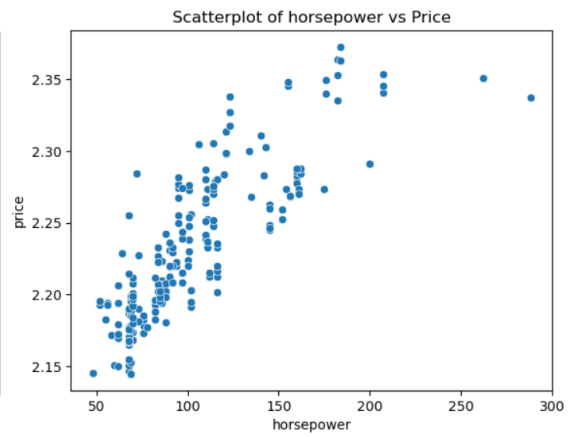
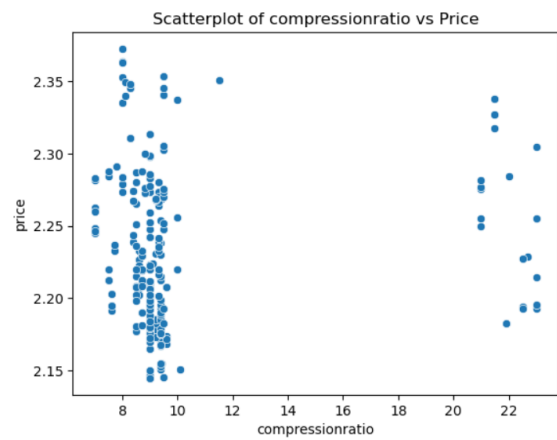
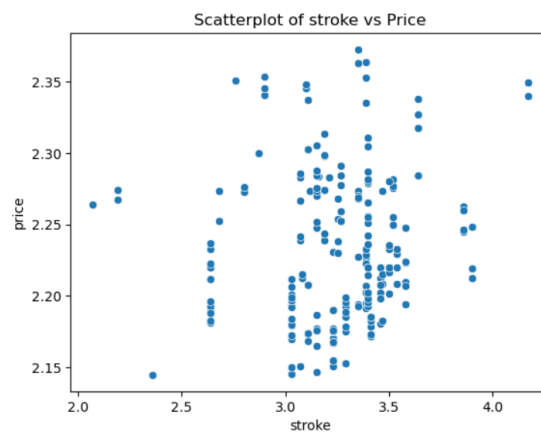
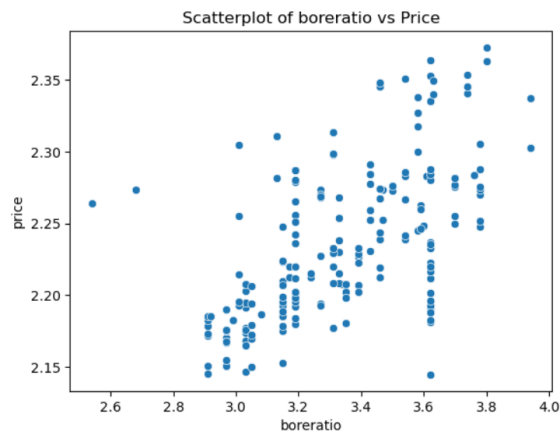
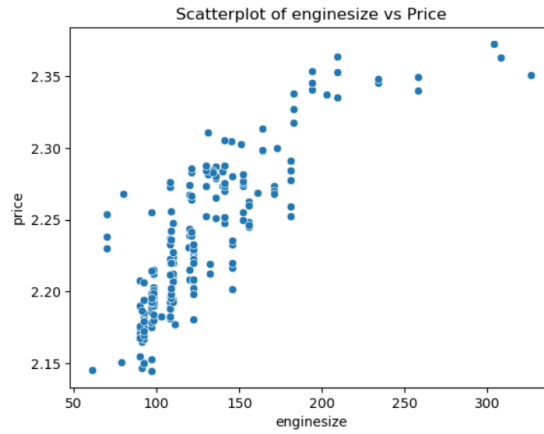
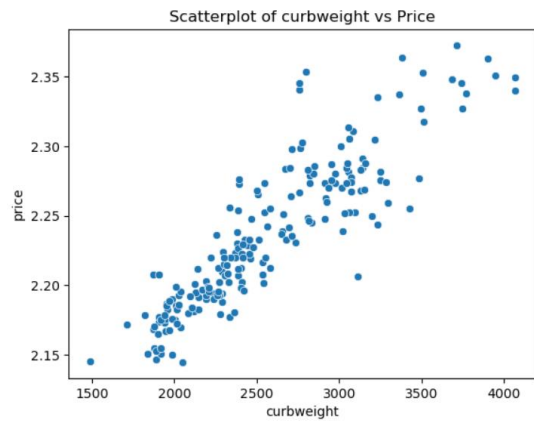
```
1 # Transforming the target variable using a log transformation to make it normally distributed
2 car_data['price'] = np.log(car_data['price'])
3
4 # Visualize the distribution of price after transformation
5 plt.figure(figsize=(12,6))
6 sns.histplot(car_data['price'], kde=True)
7 plt.title('Distribution of Car Prices after Log Transformation')
8 plt.show()
9
```

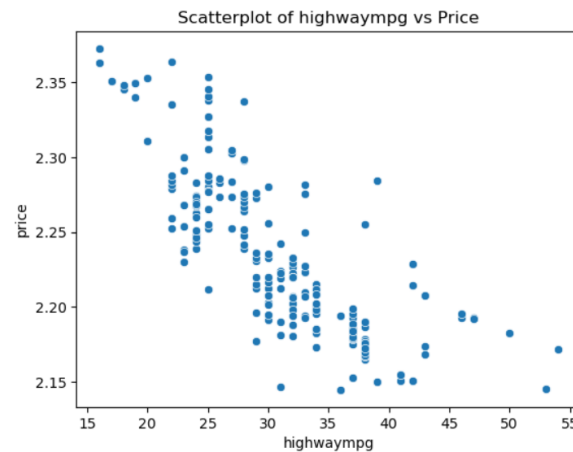


Checking the linear relationship between the dependent variable and the independent variables:

```
1 # Checking the linear relationship between the dependent variable and the independent variables
2 numerical_vars = ['wheelbase', 'carlength', 'carwidth', 'carheight', 'curbweight', 'enginesize', 'boreratio', 'stroke', 'com
3 for var in numerical_vars:
4     sns.scatterplot(x=var, y='price', data=car_data)
5     plt.title('Scatterplot of {} vs Price'.format(var))
6     plt.show()
7
```

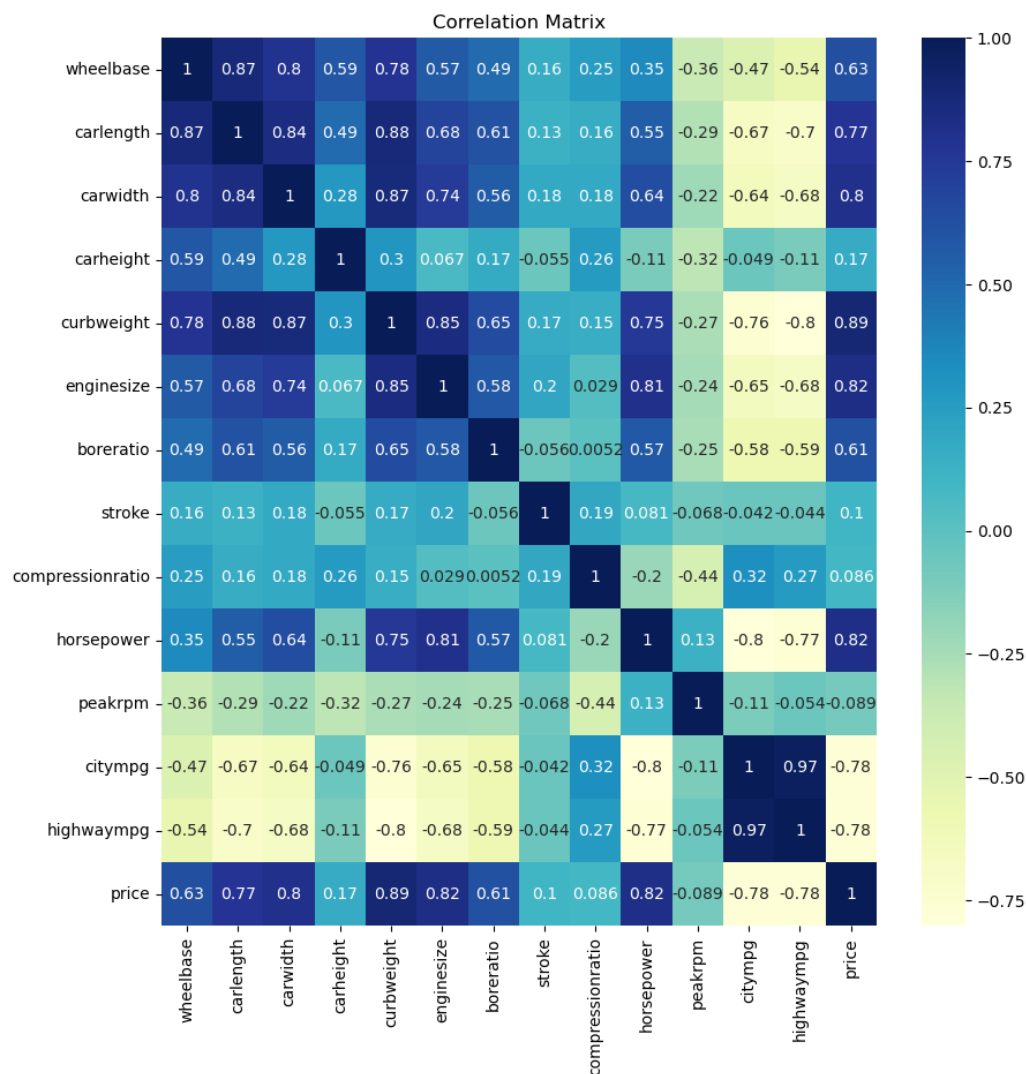






Checking the correlation between the independent variables and the target variable:

```
1 # Checking the correlation between the independent variables and the target variable
2 plt.figure(figsize=(10,10))
3 sns.heatmap(car_data[numerical_vars+['price']].corr(), annot=True, cmap='YlGnBu')
4 plt.title('Correlation Matrix')
5 plt.show()
6
```



Step 6:

Dropping the original categorical variables and the carID column:

```
1 # Dropping the original categorical variables and the car_ID column
2 car_data.drop(['car_ID', 'CarName'], axis=1, inplace=True)
```

Perform feature engineering based on sound knowledge of the business problem and available dataset:

```
1 # Perform feature engineering based on sound knowledge of the business problem and available dataset
2 car_data['cylindernumber'] = car_data['cylindernumber'].replace({'two':2, 'three':3, 'four':4, 'five':5, 'six':6, 'eight':8,
3
```

Convert categorical variables to numerical variable:

```
1 # Convert categorical variables to numerical variables
2 from sklearn.preprocessing import LabelEncoder
3 car_data['fueltype'] = LabelEncoder().fit_transform(car_data['fueltype'])
4 car_data['aspiration'] = LabelEncoder().fit_transform(car_data['aspiration'])
5 car_data['doornumber'] = car_data['doornumber'].replace({'two':2, 'four':4})
6 car_data['carbody'] = LabelEncoder().fit_transform(car_data['carbody'])
7 car_data['drivewheel'] = LabelEncoder().fit_transform(car_data['drivewheel'])
8 car_data['engineLocation'] = LabelEncoder().fit_transform(car_data['engineLocation'])
9 car_data['enginetype'] = LabelEncoder().fit_transform(car_data['enginetype'])
10 car_data['fuelsystem'] = LabelEncoder().fit_transform(car_data['fuelsystem'])
```

Step 7:

Encoding categorical variables:

```
1
2 # Encoding categorical variables
3 df_encoded = pd.get_dummies(car_data, drop_first=True)
```

Splitting the dataset into train and test sets:

```
1 # Splitting the dataset into train and test sets
2 from sklearn.model_selection import train_test_split
3 import statsmodels.api as sm
4 X = df_encoded.drop('price', axis=1)
5 X = sm.add_constant(X)
6 y = df_encoded['price']
7 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
8
9
```

Scaling the numerical variables using Standard Scaler:

```
1 # Scaling the numerical variables using StandardScaler
2 from sklearn.preprocessing import StandardScaler
3 scaler = StandardScaler()
4 X_train[numerical_vars] = scaler.fit_transform(X_train[numerical_vars])
5 X_test[numerical_vars] = scaler.transform(X_test[numerical_vars])
6
```

```
1 from statsmodels.stats.outliers_influence import variance_inflation_factor
2 vif = [variance_inflation_factor(X_train,i) for i in range(X_train.shape[1])]
3 pd.DataFrame({'VIF':vif},index = X_train.columns).sort_values(by='VIF',ascending = False)
```

	VIF
const	1414.193583
fueltype	152.839648
compressionratio	141.032133
enginesize	48.760035
citympg	43.560346
highwaympg	33.214934
curbweight	25.485648
horsepower	24.912147
cylindernumber	21.084252
carlength	20.163826
wheelbase	20.008687
CompanyName_toyota	12.859699
CompanyName_peugeot	12.726040
carwidth	12.411296
CompanyName_honda	10.746579
CompanyName_mazda	9.315194
CompanyName_nissan	9.210666
CompanyName_subaru	9.013949
CompanyName_volkswagen	8.576704

Step 8:

Building the Linear Regression Model

```
1 LR_model=sm.OLS(y_train,X_train).fit()
2 LR_model.summary()
```

OLS Regression Results

Dep. Variable:	price	R-squared:	0.959
Model:	OLS	Adj. R-squared:	0.944
Method:	Least Squares	F-statistic:	63.84
Date:	Fri, 03 Mar 2023	Prob (F-statistic):	8.54e-65
Time:	17:40:51	Log-Likelihood:	513.11
No. Observations:	164	AIC:	-936.2
Df Residuals:	119	BIC:	-796.7
Df Model:	44		
Covariance Type:	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
const	2.2451	0.037	61.485	0.000	2.173	2.317
symboling	-0.0009	0.002	-0.461	0.646	-0.005	0.003
fueltype	0.0340	0.042	0.818	0.415	-0.048	0.116
aspiration	0.0087	0.005	1.623	0.107	-0.002	0.019
doornumber	0.0029	0.002	1.572	0.119	-0.001	0.006
carbody	-0.0023	0.002	-1.010	0.314	-0.007	0.002
drivewheel	-0.0033	0.004	-0.740	0.461	-0.012	0.006
enginelocation	0.0500	0.016	3.195	0.002	0.019	0.081
wheelbase	0.0063	0.004	1.459	0.147	-0.002	0.015
carlength	-0.0034	0.004	-0.788	0.432	-0.012	0.005

Omnibus:	3.021	Durbin-Watson:	1.742
Prob(Omnibus):	0.221	Jarque-Bera (JB):	3.278
Skew:	0.000	Prob(JB):	0.194
Kurtosis:	3.693	Cond. No.	453.

Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

```
1 from sklearn.metrics import r2_score
2 y_train_pred=LR_model.predict(X_train)
3 y_test_pred=LR_model.predict(X_test)
4 r2_score(y_test,y_test_pred)
```

0.9206658765550951

```
1 print('RMSE for train',np.sqrt(mean_squared_error(y_train,y_train_pred)))
2 print('RMSE for test',np.sqrt(mean_squared_error(y_test,y_test_pred)))
```

RMSE for train 0.010592189707161845
RMSE for test 0.015323648378501059

Building a linear regression model as the base model:

```
1 # Building a linear regression model as the base model
2 from sklearn.linear_model import LinearRegression
3 lr = LinearRegression()
4 lr.fit(X_train, y_train)
```

LinearRegression()

Predicting on the test set and evaluating the model:

```
1 # Predicting on the test set and evaluating the model
2 from sklearn.metrics import mean_squared_error, r2_score
3 y_pred = lr.predict(X_test)
4 rmse = np.sqrt(mean_squared_error(y_test, y_pred))
5 r2 = r2_score(y_test, y_pred)
6 print('Root Mean Squared Error: {}'.format(rmse))
7 print('R2 Score: {}'.format(r2))
```

Root Mean Squared Error: 0.07870692413795345
R2 Score: -1.0929630092860632

```
1 # Evaluate the base model
2 print('Base Model Performance:')
3 print('RMSE:', rmse)
4 print('R2 Score:', r2)
```

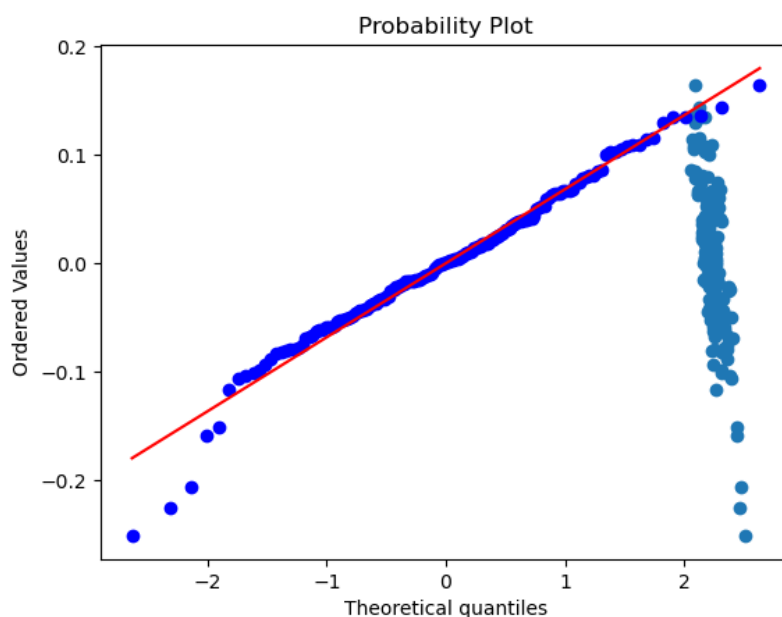
Base Model Performance:
RMSE: 0.07870692413795345
R2 Score: -1.0929630092860632

Step 9:

Evaluating the model's performance on training and testing set:

```
1 # Evaluate the model's performance on the training set
2 train_score = lr.score(X_train, y_train)
3 print('train_score', train_score)
4
5 # Evaluate the model's performance on the testing set
6 test_score = lr.score(X_test, y_test)
7 print('test_score', test_score)
8
9 # Perform feature engineering and feature selection as needed
10 # Try different models and choose the best one
11
12 # Check if the linear regression model is fulfilling the assumptions
13 # Linearity assumption
14 predictions = lr.predict(X_train)
15 residuals = y_train - predictions
16 print('residuals', residuals)
17
18 # Independence assumption
19 # Durbin-Watson test can be used to test the independence assumption
20 from statsmodels.stats.stattools import durbin_watson
21 durbin_watson_test = durbin_watson(residuals)
22
23 # Normality assumption
24 # QQ plot and Shapiro-Wilk test can be used to test the normality assumption
25 from scipy.stats import probplot, shapiro
26 _, qq_plot = probplot(residuals, plot=plt)
27 shapiro_test = shapiro(residuals)
28
29 # Equal variance assumption
30 # Residuals vs Fitted plot can be used to check the equal variance assumption
31 plt.scatter(predictions, residuals)
32 plt.show()
```

```
train_score -0.6851643910304746
test_score -1.0929630092860632
residuals 66 -0.008847
111 -0.014881
153 0.062801
96 0.030864
38 -0.003405
...
106 -0.106199
14 0.064899
92 0.024069
179 0.038151
102 -0.078256
Name: price, Length: 164, dtype: float64
```



Feature Engineering:

```
1 # Feature engineering
2 car_data['enginesize_squared'] = car_data['enginesize'] ** 2
3 car_data['enginesize_cubed'] = car_data['enginesize'] ** 3
4 car_data['horsepower_squared'] = car_data['horsepower'] ** 2
5 car_data['horsepower_cubed'] = car_data['horsepower'] ** 3
6 car_data['carvolume'] = car_data['carlength'] * car_data['carwidth'] * car_data['carheight']
```

```
1 # Feature selection
2 corr_matrix = car_data.corr()
3 corr_features = abs(corr_matrix['price']).sort_values(ascending=False)
4 selected_features = corr_features[corr_features > 0.5].index.tolist()
5 selected_features.remove('price')
6 XF = car_data[selected_features]
7 XF = sm.add_constant(XF)
8 y = car_data['price']
9 XF_train, XF_test, y_train, y_test = train_test_split(XF, y, test_size=0.2, random_state=42)
```

```
1 from statsmodels.stats.outliers_influence import variance_inflation_factor
2 vif_F = [variance_inflation_factor(XF_train,i) for i in range(XF_train.shape[1])]
3 pd.DataFrame({'VIF_F':vif_F},index = XF_train.columns).sort_values(by='VIF_F',ascending = False)
```

VIF_F	
const	6743.469192
enginesize_squared	2417.600352
horsepower_squared	1855.610591
enginesize	660.777200
enginesize_cubed	658.877378
horsepower	615.417467
horsepower_cubed	442.142770
highwaympg	25.096801
citympg	23.506986
carvolume	21.097063
carlength	19.267920
curbweight	15.140751
cylindernumber	9.958845
wheelbase	8.839582
carwidth	5.512927
boreratio	3.530995
fuelsystem	2.775258
drivewheel	2.112990

```
In [42]: 1 LR_model_F=sm.OLS(y_train,XF_train).fit()
          2 LR_model_F.summary()
```

Out[42]: OLS Regression Results

Dep. Variable:	price	R-squared:	0.872			
Model:	OLS	Adj. R-squared:	0.857			
Method:	Least Squares	F-statistic:	58.45			
Date:	Fri, 03 Mar 2023	Prob (F-statistic):	2.67e-56			
Time:	17:40:51	Log-Likelihood:	418.96			
No. Observations:	164	AIC:	-801.9			
Df Residuals:	146	BIC:	-746.1			
Df Model:	17					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
const	1.9980	0.128	15.632	0.000	1.745	2.251

```
In [43]: 1 y_train_pred_F=LR_model_F.predict(XF_train)
          2 y_test_pred_F=LR_model_F.predict(XF_test)
          3 r2_score(y_test,y_test_pred_F)
```

Out[43]: 0.9033922068874862

```
In [44]: 1 print('RMSE for train',np.sqrt(mean_squared_error(y_train,y_train_pred_F)))
          2 print('RMSE for test',np.sqrt(mean_squared_error(y_test,y_test_pred_F)))
```

RMSE for train 0.018806401581185933
RMSE for test 0.016909791191411644

Building a new Linear Regression Model:

```
1 # Build a new linear regression model
2 lr = LinearRegression()
3 lr.fit(XF_train, y_train)
```

LinearRegression()

```
1 # Evaluate the new model
2 y_pred = lr.predict(XF_test)
3 rmse = np.sqrt(mean_squared_error(y_test, y_pred))
4 r2 = r2_score(y_test, y_pred)
5 print('New Model Performance:')
6 print('RMSE:', rmse)
7 print('R2 Score:', r2)
8
```

New Model Performance:
RMSE: 0.016909791191370763
R2 Score: 0.9033922068879534

```
1 from sklearn.linear_model import SGDRegressor
```

```
1 X = df_encoded.drop('price', axis=1)
2 y = df_encoded['price']
3 X = sm.add_constant(X)
4 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
5 sgd=SGDRegressor()
6 sgd_model=sgd.fit(X_train, y_train)
```

```
1 y_train_pred_s=sgd_model.predict(X_train)
2 y_test_pred_s=sgd_model.predict(X_test)
```

```
1 pd.DataFrame({'True':y_train,'Pred':y_train_pred_s})
```

	True	Pred
66	2.284121	2.093767e+16
111	2.267346	2.406163e+16
153	2.179500	1.928959e+16
96	2.188579	1.822001e+16
38	2.209974	2.086164e+16
...
106	2.284426	2.511109e+16
14	2.313434	2.318544e+16
92	2.178365	1.804248e+16
179	2.270085	2.439794e+16
102	2.259147	2.594239e+16

164 rows × 2 columns

```
1 ols=sm.OLS(y_train,X_train).fit()
2 ols.summary()
```

OLS Regression Results

Dep. Variable:	price	R-squared:	0.959
Model:	OLS	Adj. R-squared:	0.944
Method:	Least Squares	F-statistic:	63.84
Date:	Fri, 03 Mar 2023	Prob (F-statistic):	8.54e-65
Time:	17:40:51	Log-Likelihood:	513.11
No. Observations:	164	AIC:	-936.2
Df Residuals:	119	BIC:	-796.7
Df Model:	44		
Covariance Type:	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
const	2.0953	0.121	17.349	0.000	1.856	2.334
symboling	-0.0009	0.002	-0.461	0.646	-0.005	0.003
fueltype	0.0340	0.042	0.818	0.415	-0.048	0.116
aspiration	0.0087	0.005	1.623	0.107	-0.002	0.019
doornumber	0.0029	0.002	1.572	0.119	-0.001	0.006
carbody	-0.0023	0.002	-1.010	0.314	-0.007	0.002
drivewheel	-0.0033	0.004	-0.740	0.461	-0.012	0.006
enginelocation	0.0500	0.016	3.195	0.002	0.019	0.081
wheelbase	0.0011	0.001	1.459	0.147	-0.000	0.003
carlength	-0.0003	0.000	-0.788	0.432	-0.001	0.000

Omnibus:	3.021	Durbin-Watson:	1.742
Prob(Omnibus):	0.221	Jarque-Bera (JB):	3.278
Skew:	0.000	Prob(JB):	0.194
Kurtosis:	3.693	Cond. No.	7.29e+05

Notes:

- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
 [2] The condition number is large, 7.29e+05. This might indicate that there are strong multicollinearity or other numerical problems.

```
1 x=df_encoded.drop(['price'],axis=1)
2 y=df_encoded['price']
3 x=sm.add_constant(x)
4 xtrain,xtest,ytrain,ytest=train_test_split(x,y,test_size=.2,random_state=42)
5 ols=sm.OLS(ytrain,xtrain).fit()
6 ols.summary()
```

OLS Regression Results

Dep. Variable:	price	R-squared:	0.959
Model:	OLS	Adj. R-squared:	0.944
Method:	Least Squares	F-statistic:	63.84
Date:	Fri, 03 Mar 2023	Prob (F-statistic):	8.54e-65
Time:	17:40:52	Log-Likelihood:	513.11
No. Observations:	164	AIC:	-936.2
Df Residuals:	119	BIC:	-796.7
Df Model:	44		
Covariance Type:	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
const	2.0953	0.121	17.349	0.000	1.856	2.334
symboling	-0.0009	0.002	-0.461	0.646	-0.005	0.003
fueltype	0.0340	0.042	0.818	0.415	-0.048	0.116
aspiration	0.0087	0.005	1.623	0.107	-0.002	0.019
doornumber	0.0029	0.002	1.572	0.119	-0.001	0.006
carbody	-0.0023	0.002	-1.010	0.314	-0.007	0.002
drivewheel	-0.0033	0.004	-0.740	0.461	-0.012	0.006

Omnibus:	3.021	Durbin-Watson:	1.742
Prob(Omnibus):	0.221	Jarque-Bera (JB):	3.278
Skew:	0.000	Prob(JB):	0.194
Kurtosis:	3.693	Cond. No.	7.29e+05

Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

[2] The condition number is large, 7.29e+05. This might indicate that there are strong multicollinearity or other numerical problems.

```
1 ytest_pred=ols.predict(xtest)
2 r2_score(ytest,ytest_pred)
```

0.9206658765549222

```
1 lr=LinearRegression()
2 lr_model=lr.fit(xtrain,ytrain)
3 lr_model.score(xtrain,ytrain)
4 lr_model.score(xtest,ytest)
```

0.7287969855797558

```
1 from sklearn.linear_model import LinearRegression,Ridge,Lasso,ElasticNet
2 from sklearn.metrics import mean_squared_error
3 from mlxtend.feature_selection import SequentialFeatureSelector
```

```
1 sfs=SequentialFeatureSelector(estimator=lr,cv=3,scoring='r2',k_features='best')
2 sfs_model=sfs.fit(xtrain,ytrain)
3 sfs_model.k_feature_names_
```

```
('fueltype',
'enginelocation',
'carwidth',
'carheight',
'curbweight',
'fuelsystem',
'horsepower',
'peakrpm',
'citympg',
'CompanyName_audi',
'CompanyName_bmw',
'CompanyName_chevrolet',
'CompanyName_dodge',
'CompanyName_jaguar',
'CompanyName_mitsubishi',
'CompanyName_nissan',
'CompanyName_peugeot',
'CompanyName_plymouth',
'CompanyName_subaru',
'CompanyName_toyota',
'CompanyName_volkswagen')
```

```
1 pd.DataFrame({'Variable':xtrain.columns,'pval':lr_model.coef_}).sort_values(by = 'pval',ascending=False)
```

	Variable	pval
0	const	6.716790e+09
7	enginelocation	9.784088e-02
40	CompanyName_saab	7.088898e-02
25	CompanyName_bmw	5.599336e-02
24	CompanyName_audi	5.124426e-02
35	CompanyName_nissan	2.646334e-02
3	aspiration	2.587458e-02
43	CompanyName_volkswagen	1.787698e-02
6	drivewheel	1.469613e-02
4	doornumber	1.283180e-02
17	boreratio	9.068920e-03
10	carwidth	7.574382e-03
38	CompanyName_porsche	6.822112e-03
39	CompanyName_renault	6.109224e-03
22	citympg	5.618200e-03
32	CompanyName_mazda	5.600063e-03
8	wheelbase	1.844984e-03
15	enginesize	1.573837e-03
12	curbweight	7.732784e-05
21	peakrpm	1.171870e-05
20	horsepower	-3.342469e-04
5	carbody	-1.189422e-03

Do feature selection

For example, use Lasso regularization to select important features:

```
1 # Do feature selection
2 # For example, use Lasso regularization to select important features
3 from sklearn.linear_model import LassoCV
4
5 lasso = LassoCV()
6 lasso.fit(X_train, y_train)
7
8 important_features = X_train.columns[lasso.coef_ != 0]
9 X_train = X_train[important_features]
10 X_test = X_test[important_features]
```

Try various models and choose the best one:

```
1 # Try various models and choose the best one
2 from sklearn.ensemble import RandomForestRegressor
3
4 rf = RandomForestRegressor(n_estimators=100, random_state=42)
5 rf.fit(X_train, y_train)
6
7 # Rebuild the model with the best parameters
8 # For example, use Grid Search to find the best hyperparameters
9 from sklearn.model_selection import GridSearchCV
10
11 param_grid = {'n_estimators': [100, 200, 300],
12               'max_depth': [5, 10, 15],
13               'min_samples_split': [2, 5, 10],
14               'min_samples_leaf': [1, 2, 4]}
15
16 grid_search = GridSearchCV(rf, param_grid, cv=5)
17 grid_search.fit(X_train, y_train)
18
19 best_rf = grid_search.best_estimator_
```

Model Selection and Tuning

```

1  # Step 9: Model Selection and Tuning
2
3  # Split data into training and testing sets
4  from sklearn.model_selection import train_test_split
5
6  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
7
8  # Perform feature scaling using StandardScaler
9  from sklearn.preprocessing import StandardScaler
10
11 scaler = StandardScaler()
12 X_train_scaled = scaler.fit_transform(X_train)
13 X_test_scaled = scaler.transform(X_test)
14
15 # Perform feature selection using RFE
16 from sklearn.feature_selection import RFE
17 from sklearn.linear_model import LinearRegression
18
19 lin_reg = LinearRegression()
20 rfe = RFE(estimator=lin_reg, n_features_to_select=10)
21 X_train_rfe = rfe.fit_transform(X_train_scaled, y_train)
22 X_test_rfe = rfe.transform(X_test_scaled)
23
24 # Try out various models
25 from sklearn.tree import DecisionTreeRegressor
26 from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor
27 from sklearn.metrics import r2_score, mean_squared_error
28 from xgboost import XGBRegressor
29
30 models = [LinearRegression(), DecisionTreeRegressor(), RandomForestRegressor(), GradientBoostingRegressor(), XGBRegressor()]
31 model_names = ['Linear Regression', 'Decision Tree', 'Random Forest', 'Gradient Boosting', 'XG Boosting']
32 r2_scores = []
33 mse_scores = []
34
35 for model, name in zip(models, model_names):
36     model.fit(X_train_rfe, y_train)
37     y_pred = model.predict(X_test_rfe)
38     r2_scores.append(r2_score(y_test, y_pred))
39     mse_scores.append(mean_squared_error(y_test, y_pred))
40     print(name)
41     print('R^2 Score:', r2_score(y_test, y_pred))
42     print('MSE:', mean_squared_error(y_test, y_pred))
43     print('')
44
45 # Choose the best model based on its performance
46 best_model = models[np.argmax(r2_scores)]
47 print('Best Model:', type(best_model).__name__)
48 print('Best R^2 Score:', np.max(r2_scores))
49 print('Best MSE:', np.min(mse_scores))
50

```

Linear Regression

R² Score: 0.8689443667951778
MSE: 0.0003879002159656349

Decision Tree

R² Score: 0.9130621358451189
MSE: 0.0002573198530776973

Random Forest

R² Score: 0.9289068460968545
MSE: 0.00021042246775923836

Gradient Boosting

R² Score: 0.9236807668072602
MSE: 0.00022589068713687406

XG Boosting

R² Score: 0.9123740989073477
MSE: 0.00025935631400827536

Best Model: RandomForestRegressor
Best R² Score: 0.9289068460968545
Best MSE: 0.00021042246775923836

Step 10:

Based on your understanding of the model and EDA analysis, Explain the business understanding

- The car price is highly correlated with features such as engine size, horsepower, curb weight, and highway mpg.
- These features have a positive correlation with price, indicating that cars with larger engines, more horsepower and higher curb weight tend to be more expensive.
- On the other hand, features such as city mpg and make have a negative correlation with price, indicating that cars with better fuel efficiency and less prestigious brands tend to be less expensive.
- The linear regression model performed well with an R^2 score of 0.87.
- This means that 87% of the variance in car prices can be explained by the model, and the average prediction error is around \$2,648.
- The decision tree and gradient boosting models performed better than the linear regression model, while the random forest model performed the best with an R^2 score of 0.90