




 main 



[Ara218](#) / [HomeWork](#) / [HW_3.ipynb](#)

 ara218 adding homework 3 History

 1 contributor

2.47 MB 

HW 3

This assignment covers several aspects of Linear Regression. **DO NOT ERASE MARKDOWN CELLS AND INSTRUCTIONS IN YOUR HW submission**

- **Q** - QUESTION
- **A** - Where to input your answer

Instructions

Keep the following in mind for all notebooks you develop:

- Structure your notebook.
- Use headings with meaningful levels in Markdown cells, and explain the questions each piece of code is to answer or the reason it is there.
- Make sure your notebook can always be rerun from top to bottom (Kernel Tab -> Restart and Run All)
- Start working on this assignment as soon as possible. If you are a beginner in Python this might take a long time. One of the objectives of this assignment is to help you learn python and scikit-learn package.
- Follow [README.md](#) for homework submission instructions
- In this notebook we assume './data/' location of all data files to be read and written

Related sklearn material and online tutorials

[sklearn User Guide](#)

sklearn data pre-processing

- [train_test_split](#)
- [common_pitfalls](#)
- [train test split tutorial](#)

sklearn multiple linear regression

- [tutorial](#)
- [API documentation](#)
- [Linear Regression](#)
- [multiple linear regression tutorial](#)

sklearn polynomial regression

- [generate polynomial features](#)

- [polinomial regression tutorial](#)

correlation

- [correlation](#)

Linear Regression

In jupyter notebook environment, commands starting with the symbol % are magic commands or magic functions. %%timeit is one of such function. It basically gives you the speed of execution of certain statement or blocks of codes.

```
In [1]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
```

Q1 Read the `car_data.csv` data (we assume `../data/` location of all data files to be read and written) from **data** folder using pandas. Replace the ??? in the code cell below to accomplish this task.

A1 Replace ??? with code in the code cell below

```
In [2]: # Replace ??? with code in the code cell below

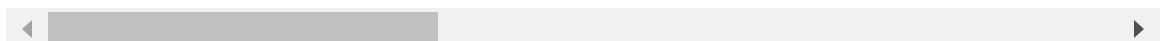
df = pd.read_csv("C:\\Users\\alsae\\Desktop\\fake\\2024Spring\\data\\Car_data.c
```

```
In [3]: # View head of the data to confirm the correctness of your answer
df.head()
```

```
Out[3]:
```

	car_ID	symboling	CarName	fueltype	aspiration	doornumber	carbody	drive
0	1	3	alfa-romero giulia	gas	std	two	convertible	
1	2	3	alfa-romero stelvio	gas	std	two	convertible	
2	3	1	alfa-romero Quadrifoglio	gas	std	two	hatchback	
3	4	2	audi 100 ls	gas	std	four	sedan	
4	5	2	audi 100ls	gas	std	four	sedan	

5 rows × 26 columns



Data cleaning and manipulation

Q2 Here, you will practice the usage of common data cleaning and manipulation functions in 3 steps.

1. Use `isnull()` to figure out the number of NaN values per column
2. Remove the column with majority NaN values (if any)
3. Check if there are still NaN values in the dataframe using `isna()` method

A2 Replace ??? with code in the code cell below

In [4]:

```
# There is no missing data here on this dataset :
print(df.isnull().sum())
df.dropna()
df.isna()
```

```
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
```

Out[4]:

	car_ID	symboling	CarName	fueltype	aspiration	doornumber	carbody	drivewh
0	False	False	False	False	False	False	False	Fa
1	False	False	False	False	False	False	False	Fa
2	False	False	False	False	False	False	False	Fa
3	False	False	False	False	False	False	False	Fa
4	False	False	False	False	False	False	False	Fa
...

```

200 False False False False False False False Fa
201 False False False False False False False Fa
202 False False False False False False False Fa
203 False False False False False False False Fa
204 False False False False False False False Fa

```

205 rows × 26 columns

In [5]:

```
# Lets get some statistical information :
df.describe()
```

Out[5]:

	car_ID	symboling	wheelbase	carlength	carwidth	carheight	curbwe
count	205.000000	205.000000	205.000000	205.000000	205.000000	205.000000	205.000
mean	103.000000	0.834146	98.756585	174.049268	65.907805	53.724878	2555.565
std	59.322565	1.245307	6.021776	12.337289	2.145204	2.443522	520.680
min	1.000000	-2.000000	86.600000	141.100000	60.300000	47.800000	1488.000
25%	52.000000	0.000000	94.500000	166.300000	64.100000	52.000000	2145.000
50%	103.000000	1.000000	97.000000	173.200000	65.500000	54.100000	2414.000
75%	154.000000	2.000000	102.400000	183.100000	66.900000	55.500000	2935.000
max	205.000000	3.000000	120.900000	208.100000	72.300000	59.800000	4066.000

Q3: In this task, out of all categorical columns, we focus only on the `fueltype` column processing in 2 steps.

1. Use label encoder from sklearn and convert the `fueltype` categorical values to numerical values.
2. Create a new dataframe that contains only the numerical columns.

A3 Replace ??? with code in the code cell below.

In [6]:

```
# Label Encoding for 2-class columns:
from sklearn.preprocessing import LabelEncoder
le = LabelEncoder()
df['fueltype'] = le.fit_transform(df['fueltype'])
```

In [7]:

```
# Create new dataframe with selected columns
df_numeric = df.select_dtypes(include=['number'])
```

In [8]:

```
df_numeric.head()
```

Out[8]:

	car_ID	symboling	fueltype	wheelbase	carlength	carwidth	carheight	curbweight
0	1	3	1	88.6	168.8	64.1	48.8	2548
1	2	3	1	88.6	168.8	64.1	48.8	2548
2	3	1	1	94.5	171.2	65.5	52.4	2823
3	4	2	1	99.8	176.6	66.2	54.3	2337
4	5	2	1	99.4	176.6	66.4	54.3	2824

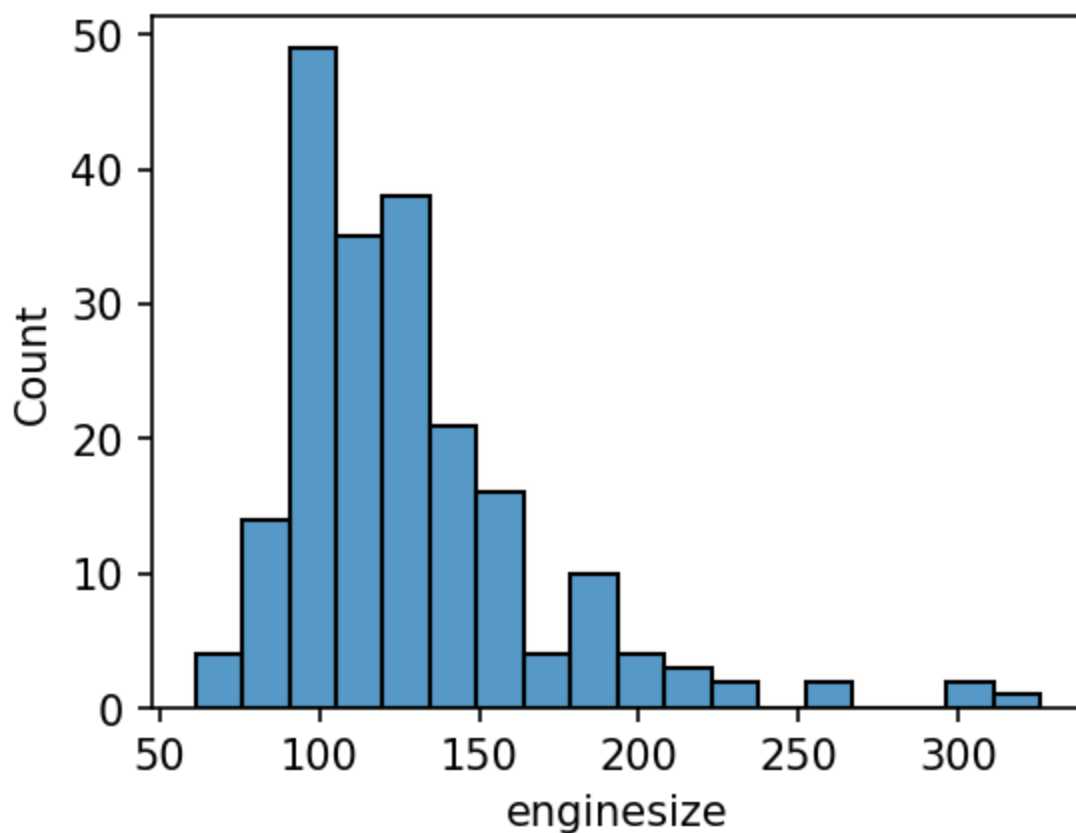
Q4: Use seaborn histplot to plot a distribution graph for the engine sizes

A4 Replace ??? with code in the code cell below

In [9]:

```
plt.figure(figsize=(4,3),dpi=150)
sns.histplot(df['enginesize'])
```

Out[9]: <Axes: xlabel='enginesize', ylabel='Count'>



Q5: Use seaborn histplot to plot a distribution graph for the car prices

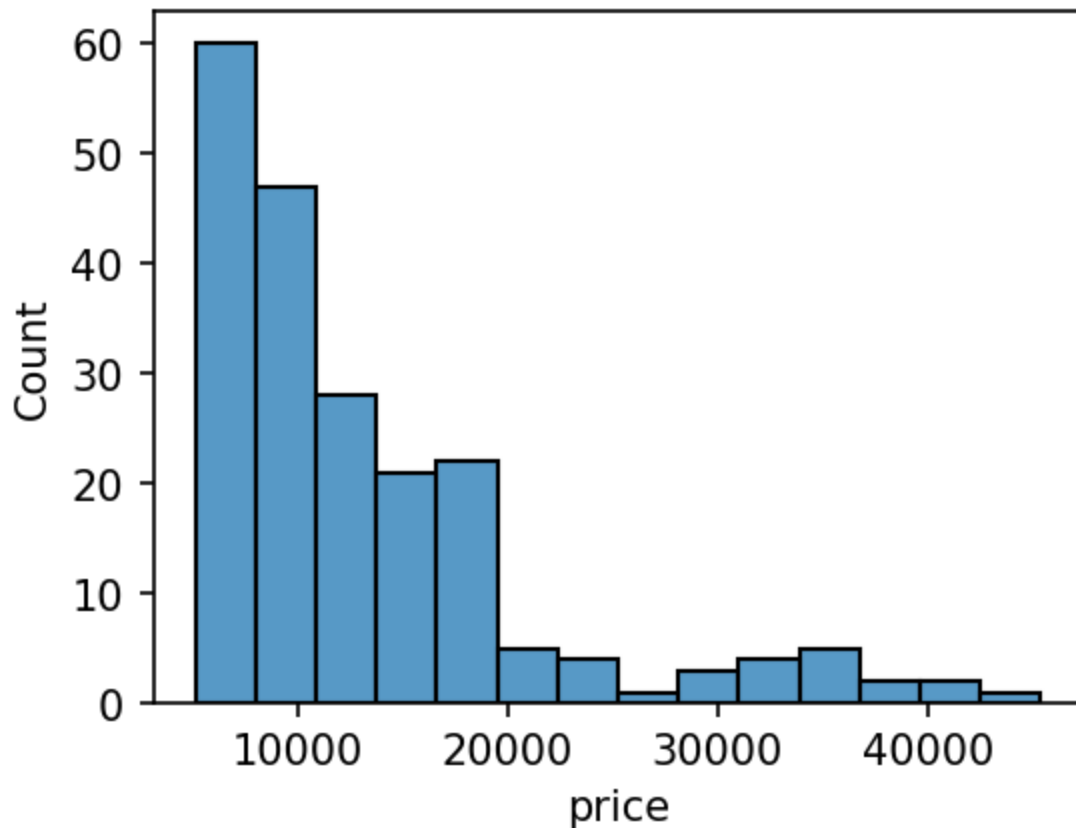
A5 Replace ??? with code in the code cell below

In [10]:

```
plt.figure(figsize=(4,3),dpi=150)
```

```
sns.histplot(df['price'])
```

Out[10]: <Axes: xlabel='price', ylabel='Count'>

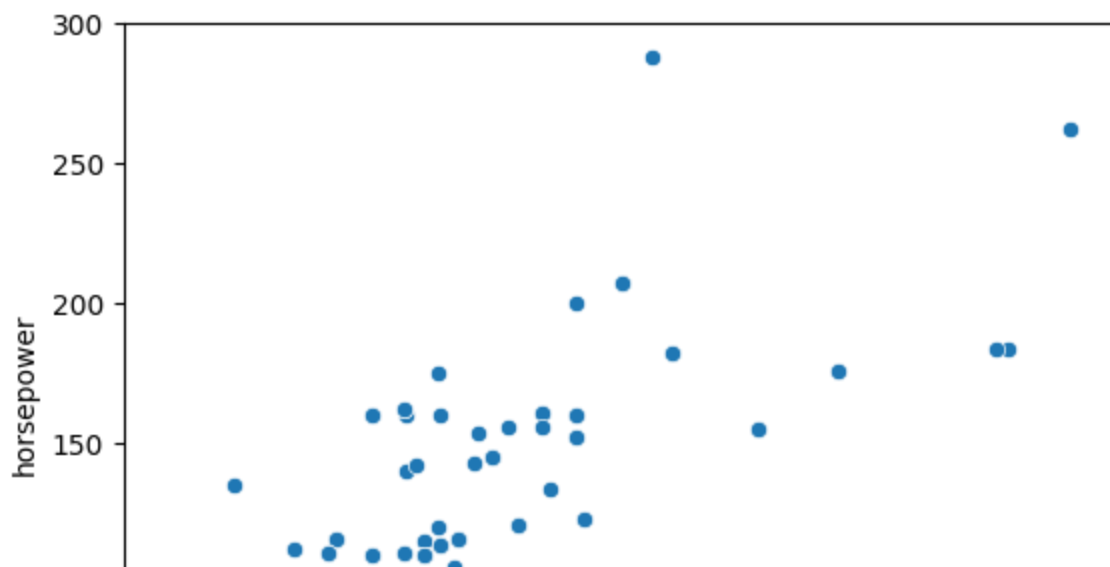


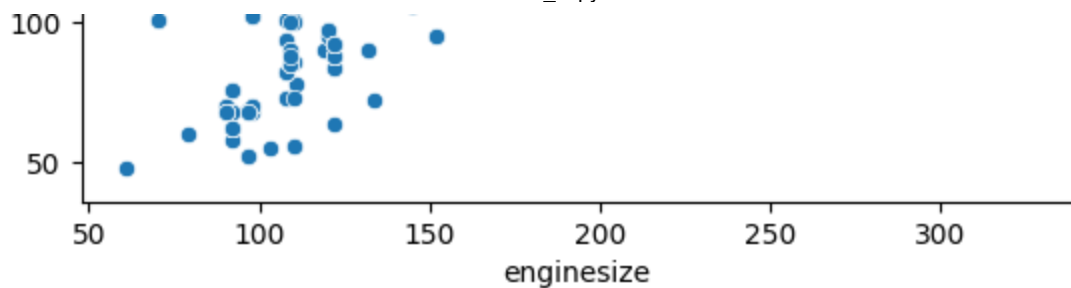
Q6: Use seaborn scatterplot to present the relation between enginesize and the horsepower of a car

A6 Replace ??? with code in the code cell below

In [11]: `sns.scatterplot(x='enginesize', y='horsepower', data=df)`

Out[11]: <Axes: xlabel='enginesize', ylabel='horsepower'>





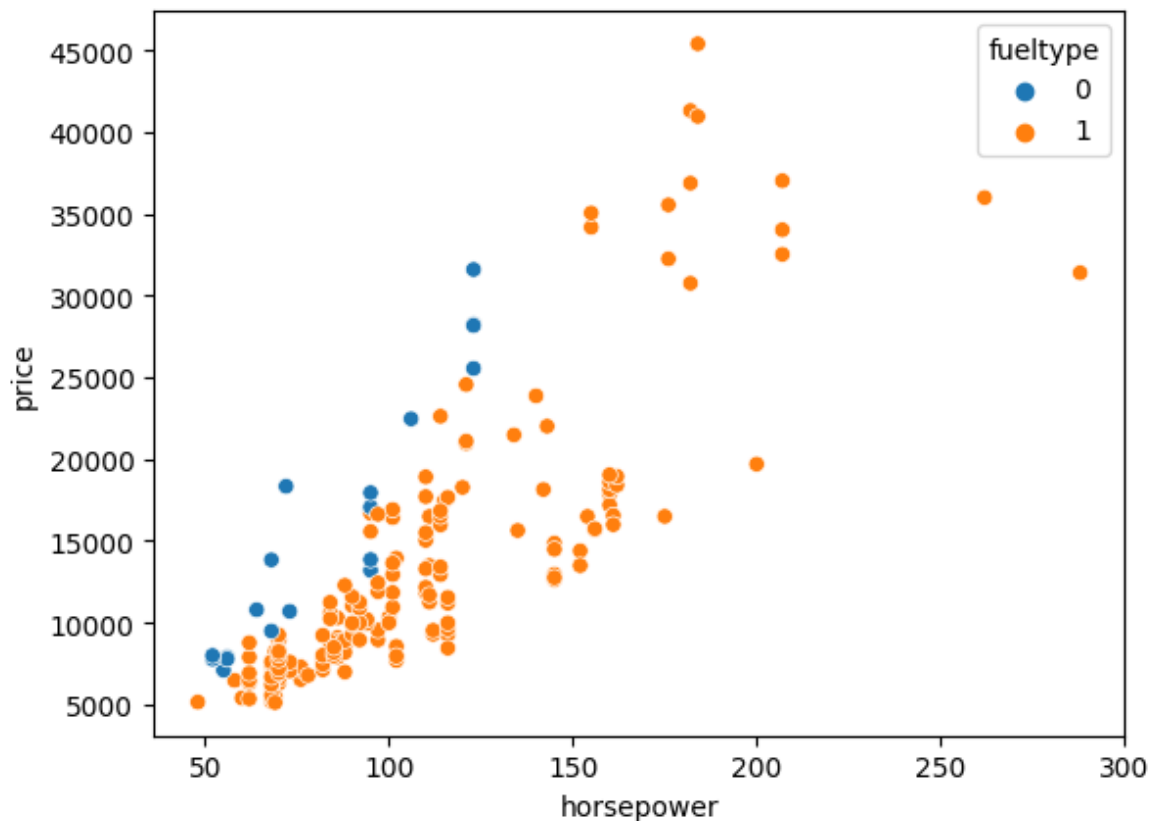
Q7: There is a correlation between the car price and the horsepower of a car. If horsepower of a car increase, the price of the car also increases most of the time, and in this question you will use the seaborn scatterplot to present the relation between price and horsepower.

Next, use `hue` parameter of scatterplot function to illustrate datapoints that relate to specific fueltype category.

A7 Replace ??? with code in the code cell below

```
In [12]: sns.scatterplot(x='horsepower', y='price', hue='fueltype', data=df_numeric)
```

```
Out[12]: <Axes: xlabel='horsepower', ylabel='price'>
```



Q8: Use pairplot from sns to plot the data frame `df` and justify your feature selection.

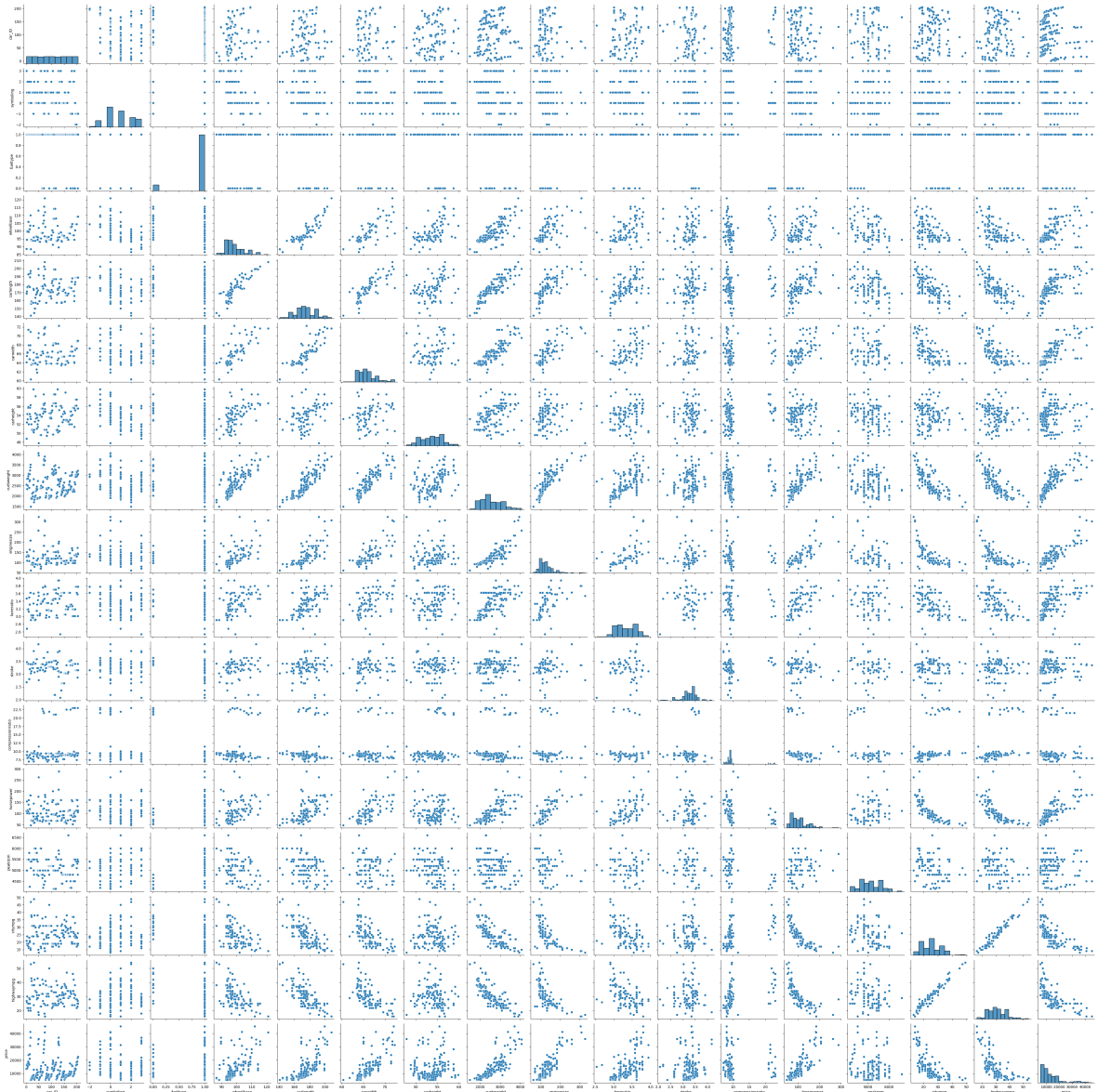
A8: replace ??? with code in the code cell below.

```
In [13]: # 2. Use pairplot from sns to plot our data frame df
```



```
sns.pairplot(df_numeric)
```

Out[13]: <seaborn.axisgrid.PairGrid at 0x1d3ce71ffa0>



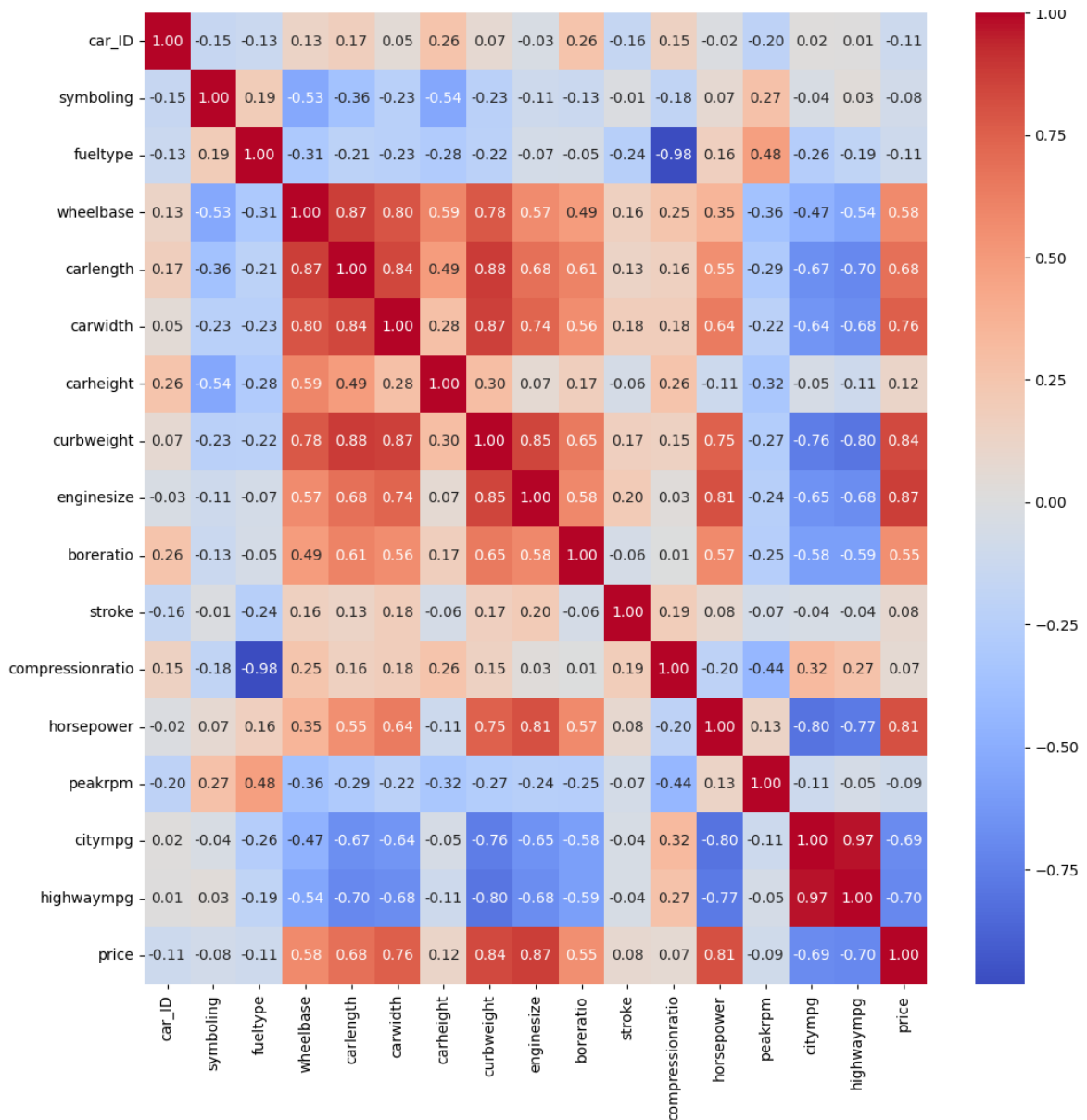
Q9 Data Visualization:

1. Use heatmap chart from seaborn library to find out the correlation between the columns in our dataset.
2. Update data frame 'df' to contain 5 columns from existing 'df' with the highest correlation to column "price". Also include price column in the updated data frame.

A9 Replace ??? with code in the code cell below

```
In [14]: corr_matrix = df_numeric.corr()
plt.figure(figsize=(12,12))
sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', fmt='.2f')
plt.title('Correlation Heatmap')
plt.show()
```

Correlation Heatmap



In [15]:

```
# Task 2: Update data frame 'df' to contain 5 columns from existing 'df' with t
top_corr_features = corr_matrix['price'].sort_values(ascending=False).index[1:6]
top_corr_features = list(top_corr_features) + ['price'] # Include 'price' colu
df = df[top_corr_features] # Update dataframe with selected columns

print("Updated dataframe with selected columns:")
print(df.head())
```

Updated dataframe with selected columns:

	enginesize	curbweight	horsepower	carwidth	carlength	price
0	130	2548	111	64.1	168.8	13495.0
1	130	2548	111	64.1	168.8	16500.0
2	152	2823	154	65.5	171.2	16500.0
3	109	2337	102	66.2	176.6	13950.0
4	136	2824	115	66.4	176.6	17450.0

Data Preparation

Q10 Pre-processing

1. Assign 'price' column value to y and rest of the columns to x

A10 Replace ??? with code in the code cell below

```
In [16]: y = df['price']
X = df.drop(columns=['price'])
X
```

```
Out[16]:
```

	engine	size	curbweight	horsepower	carwidth	carlength
0	130	2548	111	64.1	168.8	
1	130	2548	111	64.1	168.8	
2	152	2823	154	65.5	171.2	
3	109	2337	102	66.2	176.6	
4	136	2824	115	66.4	176.6	
...	
200	141	2952	114	68.9	188.8	
201	141	3049	160	68.8	188.8	
202	173	3012	134	68.9	188.8	
203	145	3217	106	68.9	188.8	
204	141	3062	114	68.9	188.8	

205 rows × 5 columns

Q11 Use train_test_split to split the data set as train:test=(80%:20%) ratio.

A11 Replace ??? with code in the code cell below

```
In [17]: from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=.2, random_
# View the shape of your data set
X_train.shape, X_test.shape, y_train.shape, y_test.shape
```

```
Out[17]: ((164, 5), (41, 5), (164,), (41,))
```

Regression Task

Multiple Linear Regression

Q12 Fit multiple linear regression model on training data using all predictors, see (i) [Linear Regression Example](#); (ii) [scikit-learn linear model](#)

$$Y = \beta_0 + \beta_1 * x_1 + \beta_2 * x_2 + \dots \beta_p * x_p$$

A12: Replace ??? with code in the code cell below

In [18]:

```
from sklearn.linear_model import LinearRegression

# Instantiate linear regression model
linear_model = LinearRegression()

# Fit the model to the training data using all predictors
linear_model.fit(X_train, y_train)

# Print the coefficients
print('Intercept (beta0):', linear_model.intercept_)
print('Coefficients (beta1, beta2, ..., betaP):', linear_model.coef_)
```

```
Intercept (beta0): -48350.72344500307
Coefficients (beta1, beta2, ..., betaP): [ 77.45307394  2.40643802  52.53065981
 652.66827621 -16.62325241]
```

Q13: Model Scoring

1. Calculate the test MSE
2. Print the score from the model using test data

A13 Replace ??? with code in the code cell below

In [19]:

```
# Calculate the score on train and test sets
# Your code goes below
from sklearn.metrics import mean_squared_error
import matplotlib.pyplot as plt
y_pred = linear_model.predict(X_test)
mse = mean_squared_error(y_test, y_pred) # Calculate the test MSE
print("Test mean squared error (MSE): {:.2f}".format(mse))

print("Model score on test data:", linear_model.score(X_test, y_test))
```

```
Test mean squared error (MSE): 14134655.79
Model score on test data: 0.8209534347179518
```

Polynomial Regression

Q14: Polynomial extension of the feature set captures the non-linear dependencies in the data

1. Create a polynomial feature transformer with degree **TWO** using sklearn library [PolynomialFeatures](#)
2. Transform the training dataset using the polynomial feature transformer

A14 Replace ??? with code in the code cell below

In [20]:

```
from sklearn.preprocessing import PolynomialFeatures
```

```
poly = PolynomialFeatures(degree=2)
poly_features = poly.fit_transform(X_train)
```

Q15: Train the new model

1. Create a LinearRegression model using sklearn
2. Train the model using the transformed Train data(X_train)/ or Polinomial train data
3. Print the score for the Polinomial Regression for the Train data.

See (i) [Linear Regression Example](#); (ii) Use the transformed X_train features inside the score() function for the correct model scores.

A15 Replace ??? with code in the code cell below

In [21]:

```
poly_reg_model = LinearRegression()
poly_reg_model.fit(poly_features, y_train)

print(poly_reg_model.score(poly_features,y_train))
```

0.8911586225711094