

# Data Analysis with Python

1. Problem statement What are the main characteristics which have the most impact on the car price ?

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

In [2]: path = "https://s3.amazonaws.com/objectstorage.softlayer.net/cf-courses-data/CognitiveClass/DA0101EN/automobileEDA.csv"
df = pd.read_csv(path)
df.head()
```

df = pd.read\_csv(path)
df.head()

Out[2]:

	symboling	normalized-losses	make	aspiration	num-of-doors	body-style	drive-wheels	engine-location	wheel-base	length	width	height	curb-weight	engine-size	bore	stroke	compression-ratio	horsepower	peak-rpm	city-mpg	highway-mpg	price	city-L100km	horsepower-binned	diesel	gas
0	3	122	alfa-romeo	std	two	convertible	rdi	front	88.6	0.81148	...	...	...	...	9.0	111.0	9000.0	21	27	1595.0	11.05476	Medium	0	1		
1	3	122	alfa-romeo	std	two	convertible	rdi	front	88.6	0.81148	...	...	...	...	9.0	111.0	9000.0	21	27	1595.0	11.05476	Medium	0	1		
2	1	122	alfa-romeo	std	two	hatchback	rdi	front	94.5	0.82281	...	...	...	...	9.0	154.0	9000.0	19	26	1650.0	12.36842	Medium	0	1		
3	2	164	audi	std	four	sedan	wdi	front	99.8	0.84603	...	...	...	...	10.0	102.0	5500.0	24	30	1395.0	9.79167	Medium	0	1		
4	2	164	audi	std	four	sedan	wdi	front	99.4	0.84603	...	...	...	...	8.0	115.0	5500.0	18	22	1745.0	13.00556	Medium	0	1		

5 rows x 29 columns

2.Analyzing Individual Feature Patterns using Visualization

In [6]: print(df.dtypes)

2. Analyzing Individual Feature Patterns using Visualization

```
In [6]: print(df.dtypes)

symboling          int64
normalized-losses  int64
make              object
aspiration        object
num-of-doors      object
drive-wheels      object
engine-location   object
wheel-base       float64
length            float64
width             float64
height            int64
curb-weight       int64
engine-type       object
num-of-cylinders  object
engine-size       int64
fuel-system       object
bore              float64
stroke            float64
compression-ratio float64
horsepower        int64
peak-rpm          float64
city-mpg          int64
highway-mpg       int64
city-L100km       float64
horsepower-binned object
diesel            int64
gas              object
dtype: object

In [7]: df.corr()
```

C:\Users\Victor\AppData\Local\Temp\ipykernel\_24488\1134722465.py:1: FutureWarning: The default value of numeric\_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric\_only to silence this warning.

	symboling	normalized-losses	wheel-base	length	width	height	curb-weight	engine-size	bore	stroke	compression-ratio	horsepower	peak-rpm	city-mpg	highway-mpg	price	city-L100km	diesel	gas
	1.000000	-0.466204	-0.535987	-0.264242	-0.500100	-0.232118	-0.110581	-0.140019	-0.092145	-0.102136	0.075819	0.279740	-0.052527	0.030223	-0.002195	0.262173	-0.196725	0.196725	
	normalized-losses	1.000000	-0.056661	0.001824	0.086802	-0.373737	0.009404	-0.112360	-0.028962	0.055563	-0.114713	0.217299	0.225943	-0.225016	-0.181877	0.032399	0.238567	-0.015448	0.015448
	wheel-base	-0.056661	1.000000	0.876024	0.814507	0.990742	0.782097	0.572027	0.489244	0.159002	0.279073	0.371147	0.360205	-0.470008	-0.543304	0.946442	0.476153	0.307727	-0.307727
	length	0.001824	0.876024	1.000000	0.857170	0.920283	0.686665	0.588025	0.489097	0.161419	0.150773	0.257612	0.268070	-0.469186	-0.686442	0.840638	0.679729	0.211197	-0.211197
	width	0.086802	0.814507	0.857170	1.000000	0.900002	0.696023	0.728436	0.544865	0.188629	0.189697	0.615077	0.245600	-0.625261	-0.686425	0.751266	0.672363	0.244356	-0.244356
	height	-0.373737	0.990742	0.920283	0.900002	1.000000	0.307381	0.074994	0.180449	-0.062704	0.259737	-0.087027	0.306974	-0.604890	-0.104812	0.325496	0.003381	0.281378	-0.281378
	curb-weight	0.009404	0.782097	0.686665	0.696023	0.307381	1.000000	0.849072	0.644060	0.167562	0.156433	0.757976	0.279740	-0.749543	-0.794889	0.834415	0.785353	0.221046	-0.221046
	engine-size	0.011236	0.572027	0.489097	0.728436	0.074994	0.849072	1.000000	0.572029	0.209622	0.001263	0.882676	0.267392	-0.605046	-0.679571	0.877235	0.745059	0.070778	-0.070778
	bore	-0.028962	0.489244	0.489097	0.544865	0.188629	0.644060	0.572809	1.000000	0.055390	0.001263	0.666636	0.267392	-0.582027	-0.591309	0.543155	0.554610	0.054648	-0.054648
	stroke	0.055563	0.159002	0.161419	0.188629	-0.062704	0.167562	0.209623	0.055390	1.000000	0.187023	0.089842	0.095713	-0.034696	-0.035201	0.082310	0.037300	0.241303	-0.241303
	compression-ratio	-0.114713	0.257612	0.268070	0.189697	0.259737	0.150773	0.000000	0.214104	-0.426780	0.314026	0.064695	0.071101	0.314026	-0.268627	0.060332	-0.060332	-0.060332	0.060332
	horsepower	0.279740	0.075819	0.217299	0.371147	0.579621	0.615077	0.679737	0.825276	0.566936	0.086462	-0.214534	1.000000	0.107985	-0.822214	-0.804675	0.898468	0.898468	-0.898468
	peak-rpm	0.279740	0.225943	0.225943	-0.307381	-0.307381	-0.307381	-0.307381	-0.307381	-0.307381	-0.307381	-0.307381	0.107985	1.000000	-0.115413	-0.059598	-0.101616	0.115413	-0.115413
	city-mpg	-0.052527	-0.225016	-0.470008	-0.685192	-0.633331	-0.649800	-0.749543	-0.650646	-0.582027	-0.034696	0.332425	-0.822214	0.115413	1.000000	0.972044	-0.688571	-0.688713	0.265676
	highway-mpg	0.030223	-0.002195	-0.118177	-0.543304	-0.686425	-0.630635	-0.784889	-0.796717	-0.591309	-0.035201	0.268465	-0.804675	-0.059598	0.972044	1.000000	-0.704692	-0.930028	0.188690
	price	-0.002195	0.032399	0.546462	0.690626	0.751266	0.135486	0.834415	0.877235	0.543155	0.052310	0.071307	0.809975	0.101616	-0.688971	-0.704692	1.000000	0.788888	0.110326
	city-L100km	0.262173	0.262173	0.262173	0.262173	0.262173	0.262173	0.262173	0.262173	0.262173	0.262173	0.262173	0.262173	0.262173	0.262173	0.262173	0.262173	0.262173	0.262173
	diesel	-0.196725	-0.015448	-0.015448	-0.015448	-0.015448	-0.015448	-0.015448	-0.015448	-0.015448	-0.015448	-0.015448	-0.015448	-0.015448	-0.015448	-0.015448	-0.015448	-0.015448	-0.015448
	gas	0.196725	0.015448	-0.015448	-0.015448	-0.015448	-0.015448	-0.015448	-0.015448	-0.015448	-0.015448	-0.015448	-0.015448	-0.015448	-0.015448	-0.015448	-0.015448	-0.015448	-0.015448

```
In [8]: # Correlation between bore, stroke, compression-ratio and horsepower.
df[['bore', 'stroke', 'compression-ratio', 'horsepower']].corr()
```

```
Out[8]:
```

	bore	stroke	compression-ratio	horsepower
bore	1.000000	-0.056390	0.001263	0.566936
stroke	-0.056390	1.000000	0.187023	0.094642
compression-ratio	0.001263	0.187023	1.000000	0.214514
horsepower	0.566936	0.094642	-0.214514	1.000000

```
In [9]: # Engine size as potential predictor variable of price
plt.figure(figsize=(8,4))
plt.plot(ymin=0)
```

```
Out[9]:
```

0.00, 53697

61366948864

50000

# Correlation between bore, stroke, compression-ratio and horsepower.  
sns.pairplot(df[['bore','stroke','compression-ratio','horsepower']], corr=True)

Out[8]:	bore	stroke	compression-ratio	horsepower
bore	1.000000	-0.003795	0.001353	0.566305
stroke	-0.003795	1.000000	0.187923	0.009462
compression-ratio	0.001353	0.187923	1.000000	-0.214514
horsepower	0.566305	0.009462	-0.214514	1.000000

# Engine size as potential predictor variable of price  
sns.regplot(x="engine-size", y="price", data=df)

Out[9]: (0.0, 53897.83265943864)



Output 9: As the engine-size goes up, the price goes up, this indicates a positive direct correlation between these two variables. Engine size seems like a pretty good predictor of price since the regression line is almost a perfect diagonal line.

# Calculating the correlation between engine and price  
df[["engine-size", "price"]].corr()

Out[10]:	engine-size	price
engine-size	1.000000	0.872335
price	0.872335	1.000000

# Highway mpg is a potential predictor variable of price  
sns.regplot(x="highway-mpg", y="price", data=df)

Out[11]: (0.0, 48171.874078142425)



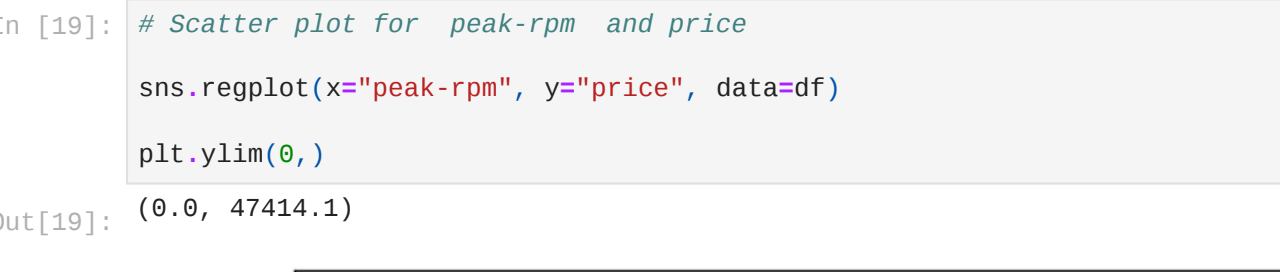
Output 11: As the highway-mpg goes up, the price goes down, this indicates an inverse/negative relationship between these two variables. Highway mpg could potentially be a predictor of price.

# Correlation between highway-mpg and price  
df[["highway-mpg", "price"]].corr()

Out[15]:	highway-mpg	price
highway-mpg	1.000000	-0.704692
price	-0.704692	1.000000

# Scatter plot for peak-rpm and price  
sns.regplot(x="peak-rpm", y="price", data=df)

Out[13]: (0.0, 47414.1)



Output 13: Peak rpm does not seem like a good predictor of the price at all since the regression line is close to horizontal. Also, the data points are very scattered and far from the fitted line, showing lots of variability. There is no a reliable variable.

# Calculating the correlation peak-rpm and price  
df[["peak-rpm", "price"]].corr()

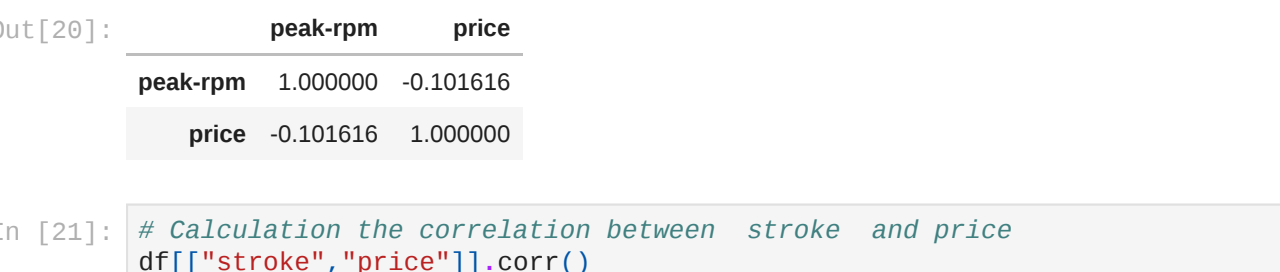
Out[14]:	peak-rpm	price
peak-rpm	1.000000	-0.101616
price	-0.101616	1.000000

# Calculating the correlation between stroke and price  
df[["stroke", "price"]].corr()

Out[21]:	stroke	price
stroke	1.000000	0.082331
price	0.082331	1.000000

sns.regplot(x="stroke", y="price", data=df)

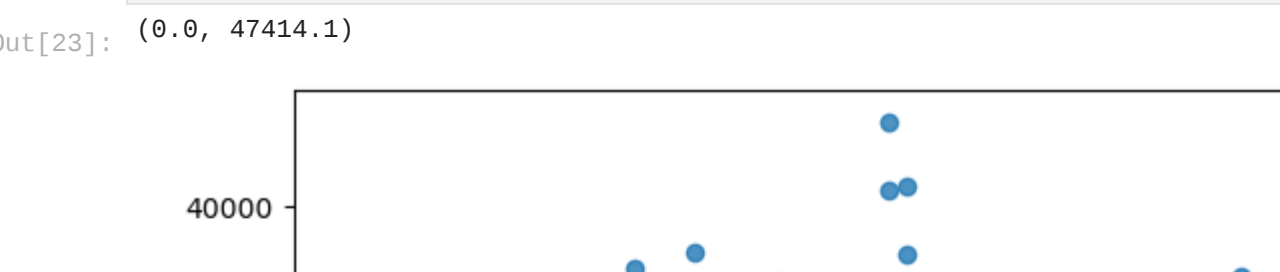
Out[22]: (0.0, 47414.1)



Output 22: Stroke is not good predictor because the linear regression line it looks more horizontal. Also scatter plots are far from the line.

# Relationship between body-style and price  
sns.boxplot(x="body-style", y="price", data=df)

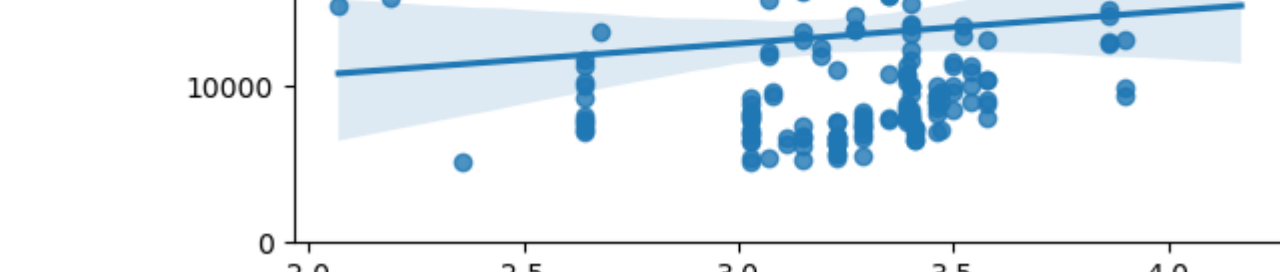
Out[7]: (0.0, 47414.1)



Output 4: It shows body\_style categories have a significant overlap, and so body-style would not be a good predictor of price.

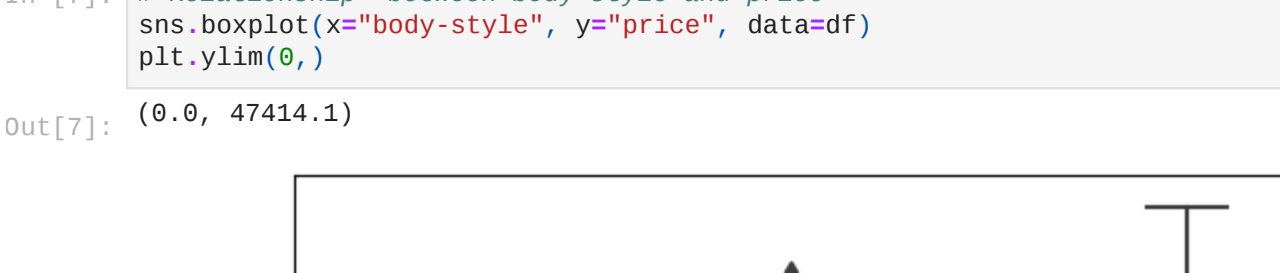
# Relationship between engine\_location and price  
sns.boxplot(x="engine-location", y="price", data=df)

Out[6]: (0.0, 47414.1)



Output 6: The distribution of price between these two engine\_location categories, front and rear, are distinct enough to take engine\_location as a potential good predictor of price.

# drive-wheels  
sns.boxplot(x="drive-wheels", y="price", data=df)



Output 8: The distribution of price between the drive-wheels categories differs as such drive-wheels could potentially be predictor of price.

price

convertible hatchback sedan wagon hardtop

body-style

body-style	min	Q1	Median	Q3	max	Outliers
convertible	12000	14000	16000	30000	31000	
hatchback	7000	8000	9000	11000	19000	21000
sedan	6000	10000	13000	18000	32000	
wagon	7000	8000	12000	16000	18000	28000
hardtop	9000	10000	20000	32000	33000	