

Table of content

1. Import Data from Module 2
2. Analyzing Individual Feature Patterns using Visualization
3. Descriptive Statistical Analysis
4. Basics of Grouping
5. Correlation and Causation
6. ANOVA

What are the main characteristics which have the most impact on the car price?

1. Import Data from Module 2

Setup
Import libraries

```
In [71]: import pandas as pd
import numpy as np
```

load data and store in dataframe df:

This dataset was hosted on IBM Cloud object click [HERE](#) for free storage

```
In [72]: path='https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBMDeveloperToolsetDemos/PyPandas/Project_Walkthrough/tutorial/P1_DataDownload/data.csv'
df = pd.read_csv(path)
df.head()
```

	symboling	normalized-losses	make	aspiration	num-of-doors	body-style	drive-wheels	engine-location	wheel-base	length	...	curb-weight	engine-size	horsepower	peak-rpm	city-mpg	highway-mpg	price
0	3	122	alfa-romero	std	two	convertible	rwd	front	88.6	0.811148	...	1489	101	113	190	24	39	13129
1	3	122	alfa-romero	std	two	convertible	rwd	front	88.6	0.811148	...	1489	101	113	190	24	39	13129
2	1	122	alfa-romero	std	two	hatchback	rwd	front	94.5	0.822681	...	1613	101	113	190	24	39	16990
3	2	164	audi	std	four	sedan	fwd	front	99.8	0.848630	...	1912	101	113	190	24	39	35900
4	2	164	audi	std	four	sedan	4wd	front	99.4	0.848630	...	1912	101	113	190	24	39	35900

5 rows × 29 columns

2. Analyzing Individual Feature Patterns using Visualization

To install seaborn we use the pip which is the python package manager.

```
In [73]: !pip install seaborn
```

Import visualization packages "Matplotlib" and "Seaborn", don't forget about "%matplotlib inline" to plot in a Jupyter notebook.

```
In [74]: import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
```

How to choose the right visualization method?

When visualizing individual variables, it is important to first understand what type of variable you are dealing with. This will help us find the right visualization method for that variable.

```
In [75]: # List the data types for each column
print(df.dtypes)
```

```
symboling          int64
normalized-losses  int64
make              object
aspiration         object
num-of-doors      object
body-style        object
drive-wheels      object
engine-location   object
wheel-base       float64
length           float64
width            float64
height           float64
curb-weight       int64
engine-size       int64
num-of-cylinders  int64
fuel-system       object
bore             float64
stroke           float64
compression-ratio float64
horsepower       float64
peak-rpm         float64
city-mpg         int64
highway-mpg      int64
price            float64
city-L/100km     float64
horsepower-binned object
diesel          int64
gas             int64
dtype: object
```

Question #1:
What is the data type of the column "peak-rpm"?

```
In [76]: # Write your code below and press Shift+Enter to execute
df['peak-rpm'].dtypes
```

Out[76]: dtype('float64')

► Click here for the solution

For example, we can calculate the correlation between variables of type "int64" or "float64" using the method "corr":

```
In [77]: df.corr()
```

	symboling	normalized-losses	wheel-base	length	width	height	curb-weight	engine-size	horsepower	peak-rpm	city-mpg	highway-mpg	price
symboling	1.000000	0.466264	-0.053597	-0.365404	-0.242423	-0.550160	-0.233118	-0.110581	-0.140019	-0.000000	0.872335	0.872335	0.872335
normalized-losses	0.466264	1.000000	-0.056661	0.019424	0.086802	-0.373737	0.099404	0.112360	0.029862	-0.000000	0.872335	0.872335	0.872335
wheel-base	-0.053597	-0.056661	1.000000	0.876024	0.814507	0.590742	0.782097	0.572027	0.4932	-0.000000	0.872335	0.872335	0.872335
length	-0.365404	0.019424	0.876024	1.000000	0.857170	0.492063	0.880665	0.685025	0.6089	-0.000000	0.872335	0.872335	0.872335
width	-0.242423	0.086802	0.814507	0.857170	1.000000	0.306002	0.866201	0.729436	0.544	-0.000000	0.872335	0.872335	0.872335
height	-0.550160	-0.373737	0.590742	0.492063	0.306002	1.000000	0.307581	0.849072	0.644	-0.000000	0.872335	0.872335	0.872335
curb-weight	-0.233118	0.099404	0.782097	0.880665	0.866201	0.307581	1.000000	0.849072	0.644	-0.000000	0.872335	0.872335	0.872335
engine-size	-0.110581	0.112360	0.572027	0.685025	0.729436	0.849072	0.849072	1.000000	0.572	-0.000000	0.872335	0.872335	0.872335
bore	-0.140019	-0.029862	0.493244	0.608971	0.544885	0.180449	0.644060	0.572609	1.000	-0.000000	0.872335	0.872335	0.872335
stroke	-0.000245	0.055563	0.158502	0.124139	0.188829	-0.062704	0.167562	0.209523	-0.0553	-0.000000	0.872335	0.872335	0.872335
compression-ratio	-0.182196	-0.114713	0.250313	0.159733	0.185967	0.259737	0.156433	0.028889	0.001	-0.000000	0.872335	0.872335	0.872335
horsepower	0.075819	0.217299	0.371147	0.579821	0.615077	-0.087027	0.757976	0.822676	0.5669	-0.000000	0.872335	0.872335	0.872335
peak-rpm	0.279740	0.239543	-0.360305	-0.285970	-0.245800	-0.309974	-0.279361	-0.256733	-0.2673	-0.000000	0.872335	0.872335	0.872335
city-mpg	-0.035527	-0.225016	-0.470606	-0.665192	-0.633531	-0.049800	-0.749543	-0.650546	-0.5820	-0.000000	0.872335	0.872335	0.872335
highway-mpg	0.062633	-0.181877	-0.543640	-0.696142	-0.680635	-0.104812	-0.794889	-0.679571	-0.5913	-0.000000	0.872335	0.872335	0.872335
price	-0.082391	0.133999	0.584642	0.690628	0.751265	0.135486	0.834415	0.872335	0.5431	-0.000000	0.872335	0.872335	0.872335
city-L/100km	0.066171	0.238567	0.476153	0.657373	0.673363	0.003811	0.785353	0.745059	0.5546	-0.000000	0.872335	0.872335	0.872335
diesel	-0.196735	-0.101546	0.307237	0.211187	0.244356	0.281578	0.221046	0.070779	0.0544	-0.000000	0.872335	0.872335	0.872335
gas	0.196735	0.101546	-0.307237	-0.211187	-0.244356	-0.281578	-0.221046	-0.070779	-0.0544	-0.000000	0.872335	0.872335	0.872335

The diagonal elements are always one, we will study correlation more precisely Pearson correlation in-depth at the end of the notebook.

Question #2:
Find the correlation between the following columns: bore, stroke, compression-ratio, and horsepower.
Hint if you would like to select those columns use the following syntax: df[['bore', 'stroke', 'compression-ratio', 'horsepower']]

```
In [78]: # Write your code below and press Shift+Enter to execute
df[['bore', 'stroke', 'compression-ratio', 'horsepower']].corr()
```

	bore	stroke	compression-ratio	horsepower
bore	1.000000	-0.055390	0.001263	0.566936
stroke	-0.055390	1.000000	0.187923	0.098462
compression-ratio	0.001263	0.187923	1.000000	-0.214514
horsepower	0.566936	0.098462	-0.214514	1.000000

► Click here for the solution

Continuous numerical variables:

Continuous numerical variables are variables that may contain any value within some range. Continuous numerical variables can have the type "int64" or "float64". A great way to visualize these variables is by using scatterplots with fitted lines.

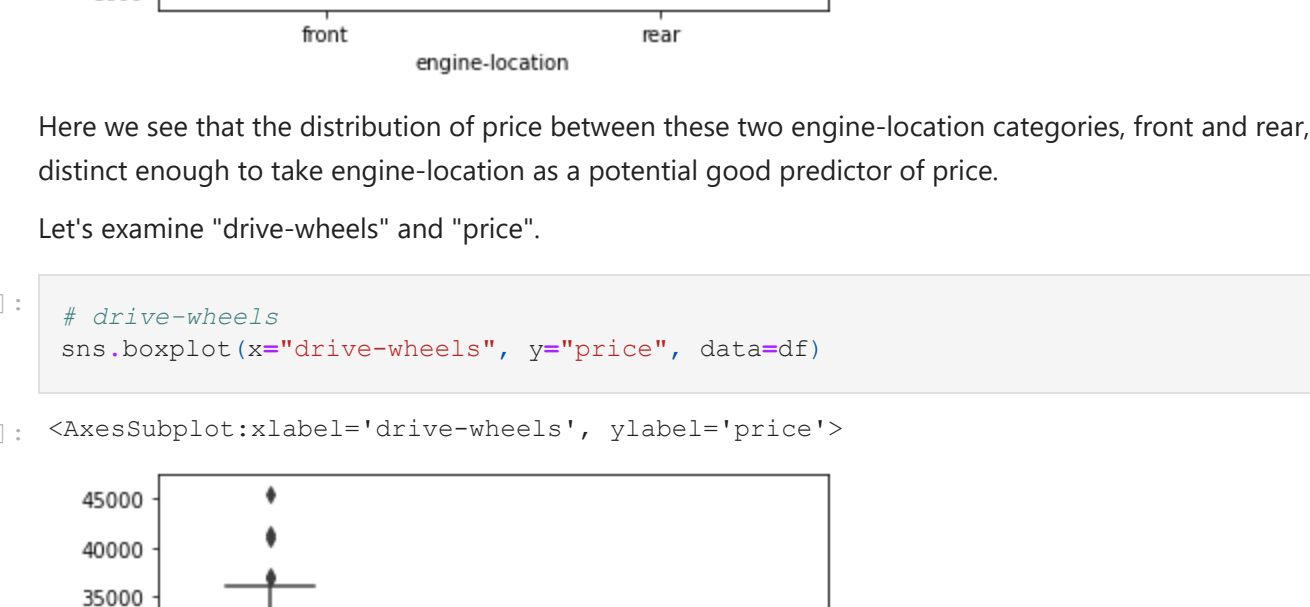
In order to start understanding the (linear) relationship between an individual variable and the price. We can do this by using "regplot", which plots the scatterplot plus the fitted regression line for the data.

Let's see several examples of different linear relationships:

Positive linear relationship

Let's find the scatterplot of "engine-size" and "price"

```
In [79]: # Engine-size as potential predictor variable of price
sns.regplot(x="engine-size", y="price", data=df)
plt.ylim(0, ...)
```



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.

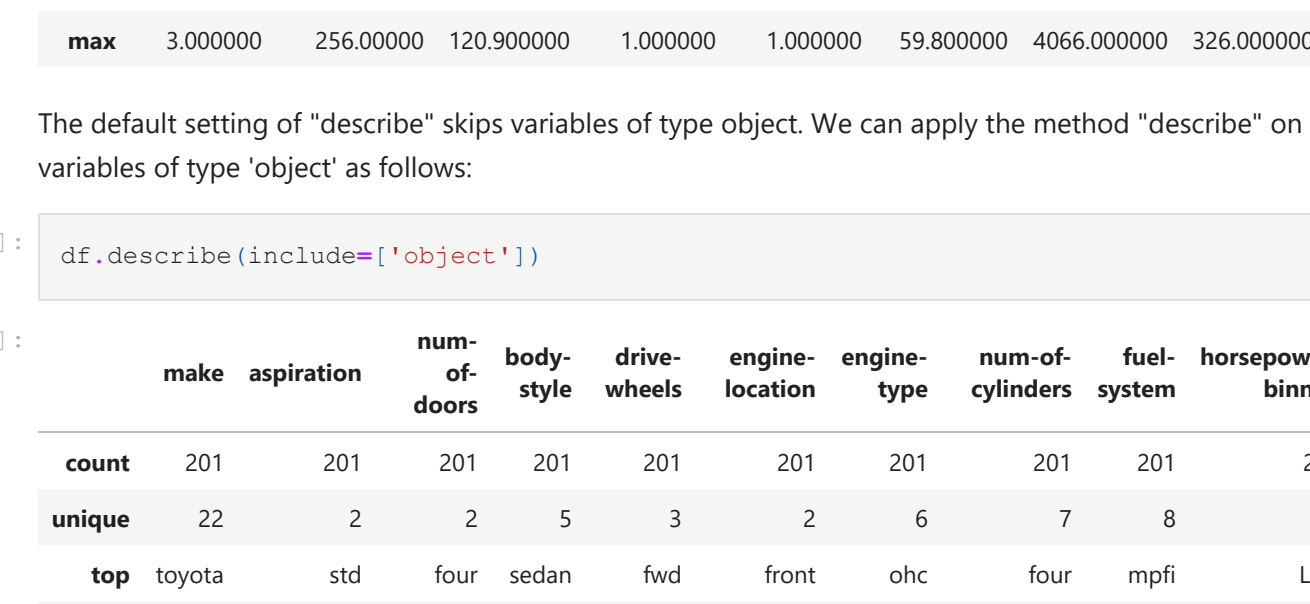
We can examine the correlation between 'engine-size' and 'price' and see it's approximately 0.87

```
In [80]: df[['engine-size', 'price']].corr()
```

	engine-size	price
engine-size	1.000000	0.872335
price	0.872335	1.000000

Highway mpg is a potential predictor variable of price

```
In [81]: sns.regplot(x="highway-mpg", y="price", data=df)
```



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.

We can examine the correlation between 'highway-mpg' and 'price' and see it's approximately -0.704

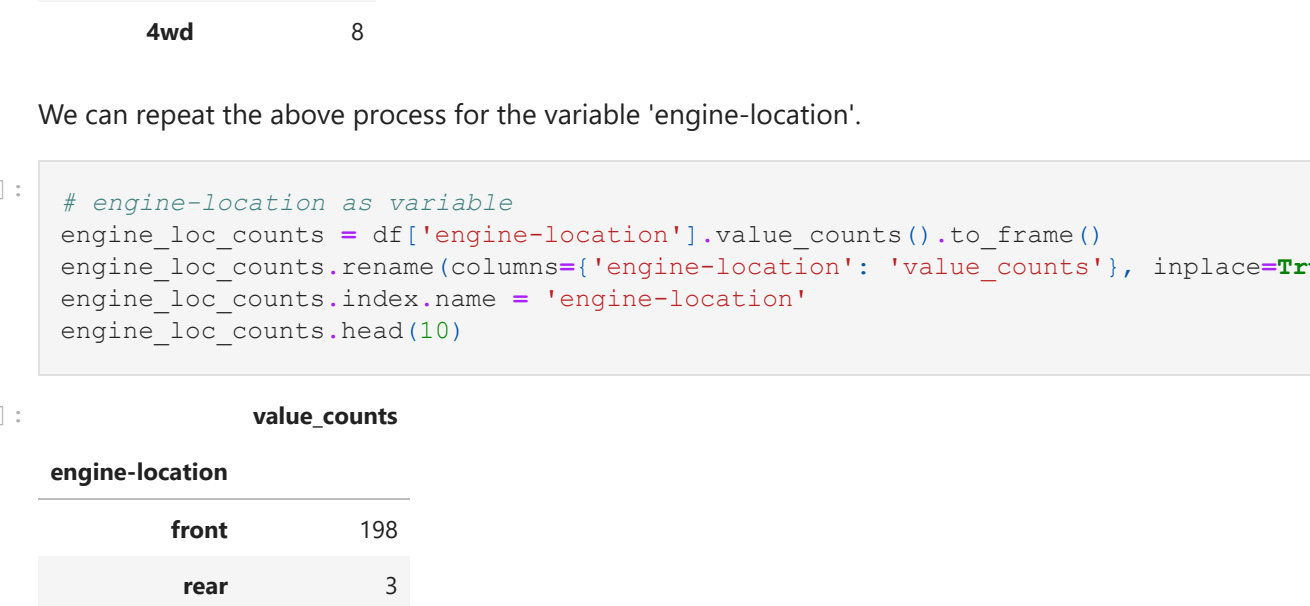
```
In [82]: df[['highway-mpg', 'price']].corr()
```

	highway-mpg	price
highway-mpg	1.000000	-0.704692
price	-0.704692	1.000000

Weak linear Relationship

Let's see if "Peak-rpm" as a predictor variable of "price".

```
In [83]: sns.regplot(x="peak-rpm", y="price", data=df)
```



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. Therefore it's not a reliable variable.

We can examine the correlation between 'peak-rpm' and 'price' and see it's approximately -0.101616

```
In [84]: df[['peak-rpm', 'price']].corr()
```

	peak-rpm	price
peak-rpm	1.000000	-0.101616
price	-0.101616	1.000000

Question 3 a):
Find the correlation between x="stroke", y="price".
Hint if you would like to select those columns use the following syntax: df[['stroke', 'price']]

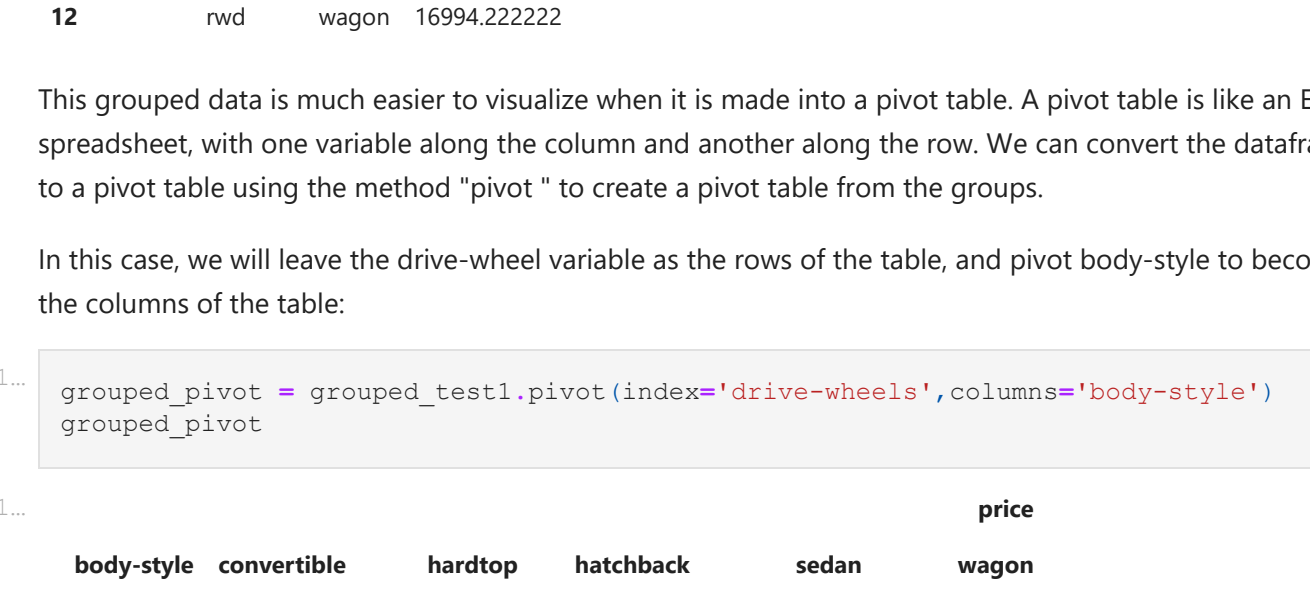
```
In [85]: # Write your code below and press Shift+Enter to execute
df[['stroke', 'price']].corr()
```

	stroke	price
stroke	1.000000	0.08231
price	0.08231	1.000000

► Click here for the solution

Question 3 b):
Given the correlation results between "price" and "stroke" do you expect a linear relationship?
Verify your results using the function "regplot()".

```
In [86]: # Write your code below and press Shift+Enter to execute
sns.regplot(x="price", y="stroke", data=df)
```



► Click here for the solution

Categorical variables

These are variables that describe a 'characteristic' of a data unit, and are selected from a small group of categories. The categorical variables can have the type "object" or "int64". A good way to visualize categorical variables is by using boxplots.

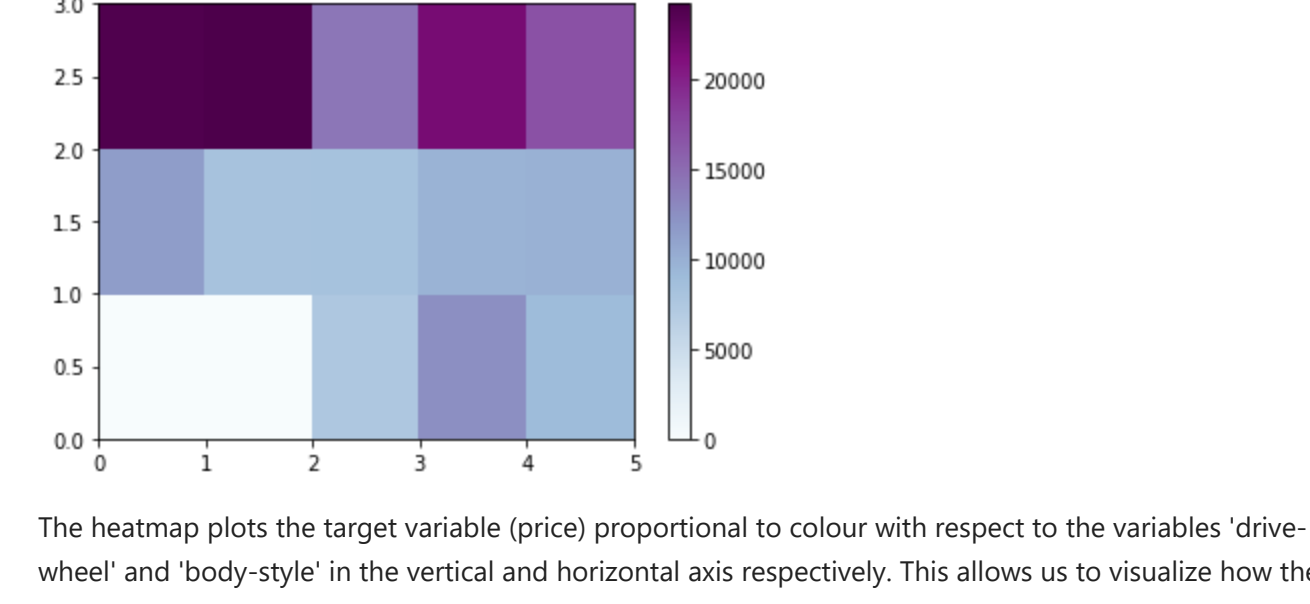
Let's look at the relationship between "body-style" and "price".

```
In [87]: sns.boxplot(x="body-style", y="price", data=df)
```



We see that the distributions of price between the different body-style categories have a significant overlap, and so body-style would not be a good predictor of price. Let's examine engine "engine-location" and "price".

```
In [88]: sns.boxplot(x="engine-location", y="price", data=df)
```



Here we see that 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.

Let's examine "drive-wheels" and "price".

```
In [89]: # drive-wheels
sns.boxplot(x="drive-wheels", y="price", data=df)
```



Here we see that the distribution of price between the different drive-wheels categories differs; as such drive-wheels could potentially be a predictor of price.

3. Descriptive Statistical Analysis

Let's first take a look at the variables by utilizing a description method.

The `describe` function automatically computes basic statistics for all continuous variables. Any NaN values are automatically skipped in these statistics.

This will show:

- the count of that variable
- the mean
- the standard deviation (std)
- the minimum value
- the IQR (Interquartile Range: 25%, 50% and 75%)
- the maximum value

We can apply the method "describe" as follows:

```
In [90]: df.describe()
```

	symboling	normalized-losses	wheel-base	length	width	height	curb-weight	engine-size
count	201.000000	201.000000	201.000000	201.000000	201.000000	201.000000	201.000000	201.000000
mean	0.840796	122.000000	98.797015	0.837102	0.915126	53.766667	2555.666667	126.875622
std	1.254802	31.95625	6.066366	0.059213	0.029187	2.447822	517.296727	41.546834
min	-2.000000	65.000000	86.600000	0.678039	0.837500	47.800000	1488.000000	61.000000
25%	0.000000	101.000000	94.500000	0.801538	0.890278	52.000000	2169.000000	98.000000
50%	1.000000	122.000000	97.000000	0.832292	0.909722	54.100000	2414.000000	120.000000
75%	2.000000	137.000000	102.400000	0.881788	0.925000	55.500000	2926.000000	141.000000
max	3.000000	256.000000	120.900000	1.000000	1.000000	59.800000	4066.000000	326.000000

The default setting of "describe" skips variables of type object. We can apply the method "describe" on the variables of type "object" as follows:

```
In [91]: df.describe(include=['object'])
```

	make	aspiration	num-of-doors	body-style	drive-wheels	engine-location	engine-type	num-of-cylinders	fuel-system	horsepower-binned
count	201	201	201	201	201	201	201	201	201	200
unique	22	2	2	5	3	2	6	7	8	3
top	toyota	std	four	sedan	fwd	front	ohc	four	mpfi	Low
freq	32	165	115	94	118	198	145	157	92	115

Value Counts

Value counts is a good way of understanding how many units of each characteristic/variable we have. We can apply the "value_counts" method on the column "drive-wheels". Don't forget the method "value_counts" only works on Pandas series, not Pandas Dataframes. As a result, we only include one bracket "df['drive-wheels']" not two brackets "df[['drive-wheels']]".

```
In [92]: df['drive-wheels'].value_counts()
```

fwd	118
rwd	75
4wd	8

Name: drive-wheels, dtype: int64

We can convert the series to a DataFrame as follows:

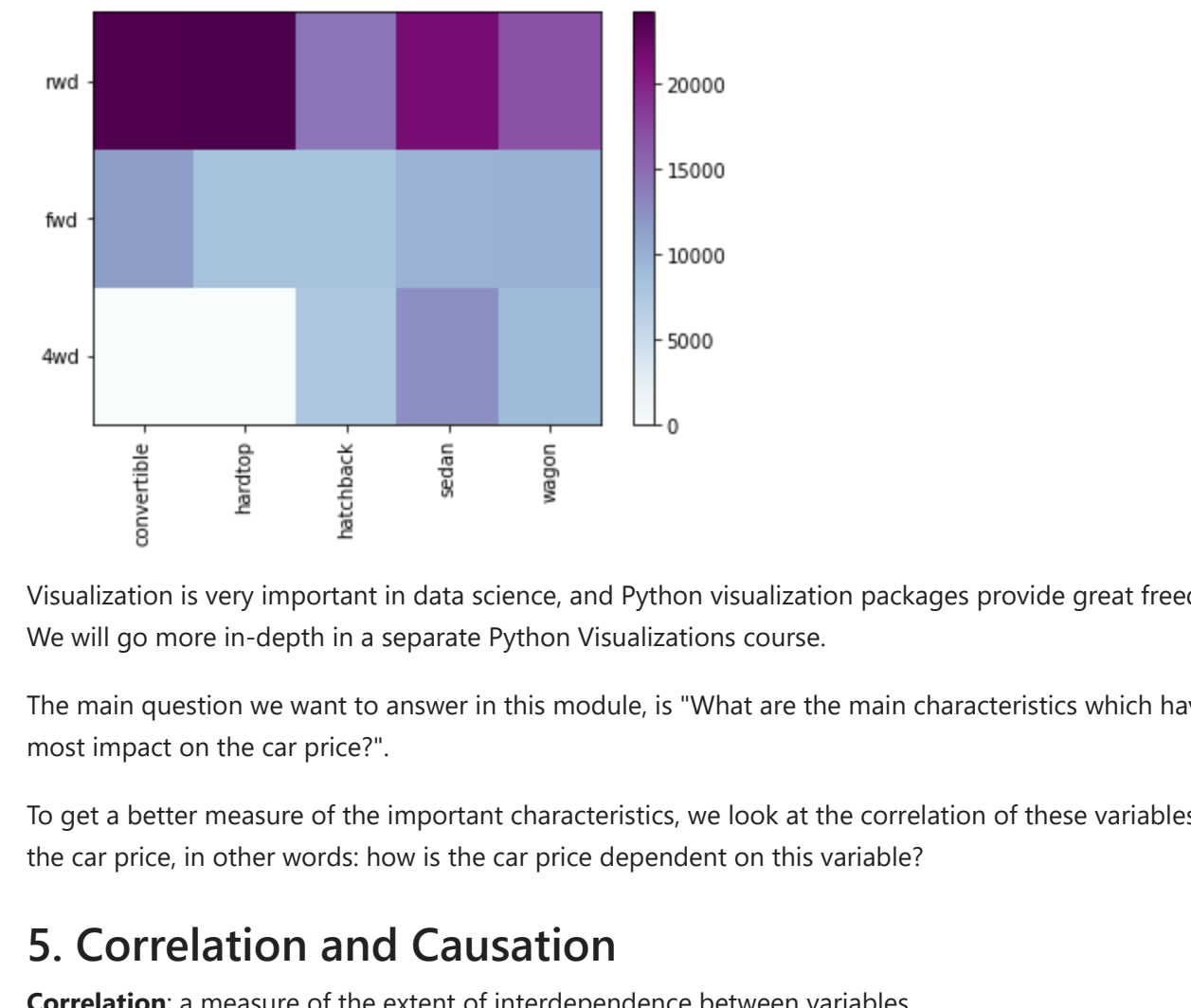
```
In [93]: df['drive-wheels'].value_counts().to_frame()
```

```
Out[93]:
```

drive-wheels	price
fwd	118
rwd	75
4wd	8

Let's repeat the above steps but save the results to the dataframe "drive_wheels_counts" and rename the column "drive-wheels" to "value_counts".

<



Visualization is very important in data science, and Python visualization packages provide great freedom. We will go more in-depth in a separate Python Visualizations course.

The main question we want to answer in this module, is "What are the main characteristics which have the most impact on the car price?".

To get a better measure of the important characteristics, we look at the correlation of these variables with the car price, in other words: how is the car price dependent on this variable?

5. Correlation and Causation

Correlation: a measure of the extent of interdependence between variables.

Causation: the relationship between cause and effect between two variables.

It is important to know the difference between these two and that correlation does not imply causation. Determining correlation is much simpler the determining causation as causation may require independent experimentation.

Pearson Correlation <p>

The Pearson Correlation measures the linear dependence between two variables X and Y.

The resulting coefficient is a value between -1 and 1 inclusive, where:

- **1:** Total positive linear correlation.
- **0:** No linear correlation, the two variables most likely do not affect each other.
- **-1:** Total negative linear correlation.

Pearson Correlation is the default method of the function 'corr'. Like before we can calculate the Pearson Correlation of the of the 'int64' or 'float64' variables.

```
In [107]: df.corr()
```

	symboling	normalized-losses	wheel-base	length	width	height	curb-weight	engine-size	bo
symboling	1.000000	0.466264	-0.535987	-0.365404	-0.242423	-0.550160	-0.233118	-0.110581	-0.1400
normalized-losses	0.466264	1.000000	-0.056661	0.019424	0.086802	-0.373737	0.099404	0.112360	-0.0298
wheel-base	-0.535987	-0.056661	1.000000	0.876024	0.814507	0.590742	0.782097	0.572027	0.4932
length	-0.365404	0.019424	0.876024	1.000000	0.857170	0.492063	0.880665	0.685025	0.6089
width	-0.242423	0.086802	0.814507	0.857170	1.000000	0.306002	0.866201	0.729436	0.5448
height	-0.550160	-0.373737	0.590742	0.492063	0.306002	1.000000	0.307581	0.074694	0.1804
curb-weight	-0.233118	0.099404	0.782097	0.880665	0.866201	0.307581	1.000000	0.849072	0.6440
engine-size	-0.110581	0.112360	0.572027	0.685025	0.729436	0.074694	0.849072	1.000000	0.5726
bore	-0.140019	-0.029862	0.493244	0.608971	0.544885	0.180449	0.644060	0.572609	1.0000
stroke	-0.008245	0.055563	0.158502	0.124139	0.188829	-0.062704	0.167562	0.209523	-0.0553
compression-ratio	-0.182196	-0.114713	0.250313	0.159733	0.189867	0.259737	0.156433	0.028889	0.0012
horsepower	0.075819	0.217299	0.371147	0.579821	0.615077	-0.087027	0.757976	0.822676	0.5669
peak-rpm	0.279740	0.239543	-0.360305	-0.285970	-0.245800	-0.309974	-0.279361	-0.256733	-0.2673
city-mpg	-0.035527	-0.225016	-0.470606	-0.665192	-0.633531	-0.049800	-0.749543	-0.650546	-0.5820
highway-mpg	0.036233	-0.181877	-0.543304	-0.698142	-0.680635	-0.104812	-0.794889	-0.679571	-0.5913
price	-0.082391	0.133999	0.584642	0.690628	0.751265	0.135486	0.834415	0.872335	0.5431
city-L/100km	0.066171	0.238567	0.476153	0.657373	0.673363	0.003811	0.785353	0.745059	0.5546
diesel	-0.196735	-0.101546	0.307237	0.211187	0.244356	0.281578	0.221046	0.070779	0.0544
gas	0.196735	0.101546	-0.307237	-0.211187	-0.244356	-0.281578	-0.221046	-0.070779	-0.0544

sometimes we would like to know the significant of the correlation estimate.

P-value:

What is this P-value? The P-value is the probability value that the correlation between these two variables is statistically significant. Normally, we choose a significance level of 0.05, which means that we are 95% confident that the correlation between the variables is significant.

By convention, when the

- p-value is < 0.001: we say there is strong evidence that the correlation is significant.
- the p-value is < 0.05: there is moderate evidence that the correlation is significant.
- the p-value is < 0.1: there is weak evidence that the correlation is significant.
- the p-value is > 0.1: there is no evidence that the correlation is significant.

We can obtain this information using 'stats' module in the 'scipy' library.

```
In [108]: from scipy import stats
```

	symboling	normalized-losses	wheel-base	length	width	height	curb-weight	engine-size	bo
symboling	1.000000	0.466264	-0.535987	-0.365404	-0.242423	-0.550160	-0.233118	-0.110581	-0.1400
normalized-losses	0.466264	1.000000	-0.056661	0.019424	0.086802	-0.373737	0.099404	0.112360	-0.0298
wheel-base	-0.535987	-0.056661	1.000000	0.876024	0.814507	0.590742	0.782097	0.572027	0.4932
length	-0.365404	0.019424	0.876024	1.000000	0.857170	0.492063	0.880665	0.685025	0.6089
width	-0.242423	0.086802	0.814507	0.857170	1.000000	0.306002	0.866201	0.729436	0.5448
height	-0.550160	-0.373737	0.590742	0.492063	0.306002	1.000000	0.307581	0.074694	0.1804
curb-weight	-0.233118	0.099404	0.782097	0.880665	0.866201	0.307581	1.000000	0.849072	0.6440
engine-size	-0.110581	0.112360	0.572027	0.685025	0.729436	0.074694	0.849072	1.000000	0.5726
bore	-0.140019	-0.029862	0.493244	0.608971	0.544885	0.180449	0.644060	0.572609	1.0000
stroke	-0.008245	0.055563	0.158502	0.124139	0.188829	-0.062704	0.167562	0.209523	-0.0553
compression-ratio	-0.182196	-0.114713	0.250313	0.159733	0.189867	0.259737	0.156433	0.028889	0.0012
horsepower	0.075819	0.217299	0.371147	0.579821	0.615077	-0.087027	0.757976	0.822676	0.5669
peak-rpm	0.279740	0.239543	-0.360305	-0.285970	-0.245800	-0.309974	-0.279361	-0.256733	-0.2673
city-mpg	-0.035527	-0.225016	-0.470606	-0.665192	-0.633531	-0.049800	-0.749543	-0.650546	-0.5820
highway-mpg	0.036233	-0.181877	-0.543304	-0.698142	-0.680635	-0.104812	-0.794889	-0.679571	-0.5913
price	-0.082391	0.133999	0.584642	0.690628	0.751265	0.135486	0.834415	0.872335	0.5431
city-L/100km	0.066171	0.238567	0.476153	0.657373	0.673363	0.003811	0.785353	0.745059	0.5546
diesel	-0.196735	-0.101546	0.307237	0.211187	0.244356	0.281578	0.221046	0.070779	0.0544
gas	0.196735	0.101546	-0.307237	-0.211187	-0.244356	-0.281578	-0.221046	-0.070779	-0.0544

Wheel-base vs Price

Let's calculate the Pearson Correlation Coefficient and P-value of 'wheel-base' and 'price'.

```
In [109]: pearson_coef, p_value = stats.pearsonr(df['wheel-base'], df['price'])
print("The Pearson Correlation Coefficient is", pearson_coef, " with a P-value of P = "
```

	symboling	normalized-losses	wheel-base	length	width	height	curb-weight	engine-size	bo
symboling	1.000000	0.466264	-0.535987	-0.365404	-0.242423	-0.550160	-0.233118	-0.110581	-0.1400
normalized-losses	0.466264	1.000000	-0.056661	0.019424	0.086802	-0.373737	0.099404	0.112360	-0.0298
wheel-base	-0.535987	-0.056661	1.000000	0.876024	0.814507	0.590742	0.782097	0.572027	0.4932
length	-0.365404	0.019424	0.876024	1.000000	0.857170	0.492063	0.880665	0.685025	0.6089
width	-0.242423	0.086802	0.814507	0.857170	1.000000	0.306002	0.866201	0.729436	0.5448
height	-0.550160	-0.373737	0.590742	0.492063	0.306002	1.000000	0.307581	0.074694	0.1804
curb-weight	-0.233118	0.099404	0.782097	0.880665	0.866201	0.307581	1.000000	0.849072	0.6440
engine-size	-0.110581	0.112360	0.572027	0.685025	0.729436	0.074694	0.849072	1.000000	0.5726
bore	-0.140019	-0.029862	0.493244	0.608971	0.544885	0.180449	0.644060	0.572609	1.0000
stroke	-0.008245	0.055563	0.158502	0.124139	0.188829	-0.062704	0.167562	0.209523	-0.0553
compression-ratio	-0.182196	-0.114713	0.250313	0.159733	0.189867	0.259737	0.156433	0.028889	0.0012
horsepower	0.075819	0.217299	0.371147	0.579821	0.615077	-0.087027	0.757976	0.822676	0.5669
peak-rpm	0.279740	0.239543	-0.360305	-0.285970	-0.245800	-0.309974	-0.279361	-0.256733	-0.2673
city-mpg	-0.035527	-0.225016	-0.470606	-0.665192	-0.633531	-0.049800	-0.749543	-0.650546	-0.5820
highway-mpg	0.036233	-0.181877	-0.543304	-0.698142	-0.680635	-0.104812	-0.794889	-0.679571	-0.5913
price	-0.082391	0.133999	0.584642	0.690628	0.751265	0.135486	0.834415	0.872335	0.5431
city-L/100km	0.066171	0.238567	0.476153	0.657373	0.673363	0.003811	0.785353	0.745059	0.5546
diesel	-0.196735	-0.101546	0.307237	0.211187	0.244356	0.281578	0.221046	0.070779	0.0544
gas	0.196735	0.101546	-0.307237	-0.211187	-0.244356	-0.281578	-0.221046	-0.070779	-0.0544

Conclusion:

Since the p-value is < 0.001, the correlation between wheel-base and price is statistically significant, although the linear relationship isn't extremely strong (~0.585)

Horsepower vs Price

Let's calculate the Pearson Correlation Coefficient and P-value of 'horsepower' and 'price'.

```
In [110]: pearson_coef, p_value = stats.pearsonr(df['horsepower'], df['price'])
print("The Pearson Correlation Coefficient is", pearson_coef, " with a P-value of P = "
```

	symboling	normalized-losses	wheel-base	length	width	height	curb-weight	engine-size	bo
symboling	1.000000	0.466264	-0.535987	-0.365404	-0.242423	-0.550160	-0.233118	-0.110581	-0.1400
normalized-losses	0.466264	1.000000	-0.056661	0.019424	0.086802	-0.373737	0.099404	0.112360	-0.0298
wheel-base	-0.535987	-0.056661	1.000000	0.876024	0.814507	0.590742	0.782097	0.572027	0.4932
length	-0.365404	0.019424	0.876024	1.000000	0.857170	0.492063	0.880665	0.685025	0.6089
width	-0.242423	0.086802	0.814507	0.857170	1.000000	0.306002	0.866201	0.729436	0.5448
height	-0.550160	-0.373737	0.590742	0.492063	0.306002	1.000000	0.307581	0.074694	0.1804
curb-weight	-0.233118	0.099404	0.782097	0.880665	0.866201	0.307581	1.000000	0.849072	0.6440
engine-size	-0.110581	0.112360	0.572027	0.685025	0.729436	0.074694	0.849072	1.000000	0.5726
bore	-0.140019	-0.029862	0.493244	0.608971	0.544885	0.180449	0.644060	0.572609	1.0000
stroke	-0.008245	0.055563	0.158502	0.124139	0.188829	-0.062704	0.167562	0.209523	-0.0553
compression-ratio	-0.182196	-0.114713	0.250313	0.159733	0.189867	0.259737	0.156433	0.028889	0.0012
horsepower	0.075819	0.217299	0.371147	0.579821	0.615077	-0.087027	0.757976	0.822676	0.5669
peak-rpm	0.279740	0.239543	-0.360305	-0.285970	-0.245800	-0.309974	-0.279361	-0.256733	-0.2673
city-mpg	-0.035527	-0.225016	-0.470606	-0.665192	-0.633531	-0.049800	-0.749543	-0.650546	-0.5820
highway-mpg	0.036233	-0.181877	-0.543304	-0.698142	-0.680635	-0.104812	-0.794889	-0.679571	-0.5913
price	-0.082391	0.133999	0.584642	0.690628	0.751265	0.135486	0.834415	0.872335	0.5431
city-L/100km	0.066171	0.238567	0.476153	0.657373	0.673363	0.003811	0.785353	0.745059	0.5546
diesel	-0.196735	-0.101546	0.307237	0.211187	0.244356	0.281578	0.221046	0.070779	0.0544
gas	0.196735	0.101546	-0.307237	-0.211187	-0.244356	-0.281578	-0.221046	-0.070779	-0.0544

Conclusion:

Since the p-value is < 0.001, the correlation between horsepower and price is statistically significant, and the linear relationship is quite strong (~0.809, close to 1)

Length vs Price

Let's calculate the Pearson Correlation Coefficient and P-value of 'length' and 'price'.

```
In [111]: pearson_coef, p_value = stats.pearsonr(df['length'], df['price'])
print("The Pearson Correlation Coefficient is", pearson_coef, " with a P-value of P = "
```

	symboling	normalized-losses	wheel-base	length	width	height	curb-weight	engine-size	bo
symboling	1.000000	0.466264	-0.535987	-0.365404	-0.242423	-0.550160	-0.233118	-0.110581	-0.1400
normalized-losses	0.466264	1.000000	-0.056661	0.019424	0.086802	-0.373737	0.099404	0.112360	-0.0298