

# Exp-7: Study of various useful Exploratory Data Analysis (EDA) techniques in Python using Pandas, NumPy, Matplotlib and Seaborn.

## Exploratory Data Analysis (EDA):

EDA is **preliminary step** in data analysis in-order to:

1. summarize main characteristics of the data
2. gain better understanding of the data set
3. uncover relationship between different variables and
4. extract important variables for the problem to solve.

## Objectives:

1. Descriptive statistics
  - Describe basic features of data
  - Giving short summaries about the sample and measures of data
2. Groupby
3. Correlation (basic)
4. Correlation- Statistics (Pearson correlation and correlation heatmaps) Explore features or characteristics to predict price of car

## Table of Contents

1. [Import Data from Module](#)
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 that have the most impact on the car price?

# 1. Import Data from Module 2

## Setup

Import libraries:

```
#install specific version of libraries used in lab
#! mamba install pandas==1.3.3
#! mamba install numpy=1.21.2
#! mamba install scipy=1.7.1-y
#! mamba install seaborn=0.9.0-y
```

```
import pandas as pd
import numpy as np
```

Load the data and store it in dataframe df:

This dataset is the output csv of Experiment 5 that we had performed.

```
#path=                                         % Specify the path of the dataset 'automobileEDA.csv'
df = pd.read_csv('/content/automobileEDA.csv')
df.head()
```

	symboling	normalized-losses	make	aspiration	num-of-doors	body-style	drive-wheels	engine-location	wheel-base
0	3	122	alfa-romero		std	two	convertible	rwd	front
1	3	122	alfa-romero		std	two	convertible	rwd	front
2	1	122	alfa-romero		std	two	hatchback	rwd	front
3	2	164	audi		std	four	sedan	fwd	front
4	2	164	audi		std	four	sedan	4wd	front

5 rows × 29 columns



## 2. Analyzing Individual Feature Patterns Using Visualization

To install Seaborn we use pip, the Python package manager.

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

```
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.

```
# 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-type	object
num-of-cylinders	object
engine-size	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"?**

```
# Write your code below and press Shift+Enter to execute
print(df.dtypes['peak-rpm'])
```

```
float64
```

► Click here for the result

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

```
df.corr()
```

	<b>symboling</b>	<b>normalized-losses</b>	<b>wheel-base</b>	<b>length</b>	<b>width</b>	<b>height</b>	<b>curb-weight</b>
<b>symboling</b>	1.000000	0.466264	-0.535987	-0.365404	-0.242423	-0.550160	-0.233118
<b>normalized-losses</b>	0.466264	1.000000	-0.056661	0.019424	0.086802	-0.373737	0.099404
<b>wheel-base</b>	-0.535987	-0.056661	1.000000	0.876024	0.814507	0.590742	0.782097
<b>length</b>	-0.365404	0.019424	0.876024	1.000000	0.857170	0.492063	0.880665
<b>width</b>	-0.242423	0.086802	0.814507	0.857170	1.000000	0.306002	0.866201
<b>height</b>	-0.550160	-0.373737	0.590742	0.492063	0.306002	1.000000	0.307581

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']]
```

```
cit魯mpg      0.035597    0.225016    0.170606    0.665102    0.622521    0.010200    0.710515
# Write your code below and press Shift+Enter to execute
df[['bore','stroke','compression-ratio','horsepower']].corr()
```

	<b>bore</b>	<b>stroke</b>	<b>compression-ratio</b>	<b>horsepower</b>	
<b>bore</b>	1.000000	-0.055390	0.001263	0.566936	
<b>stroke</b>	-0.055390	1.000000	0.187923	0.098462	
<b>compression-ratio</b>	0.001263	0.187923	1.000000	-0.214514	
<b>horsepower</b>	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. They can be of type "int64" or "float64". A great way to visualize these variables is by using scatterplots with fitted lines.

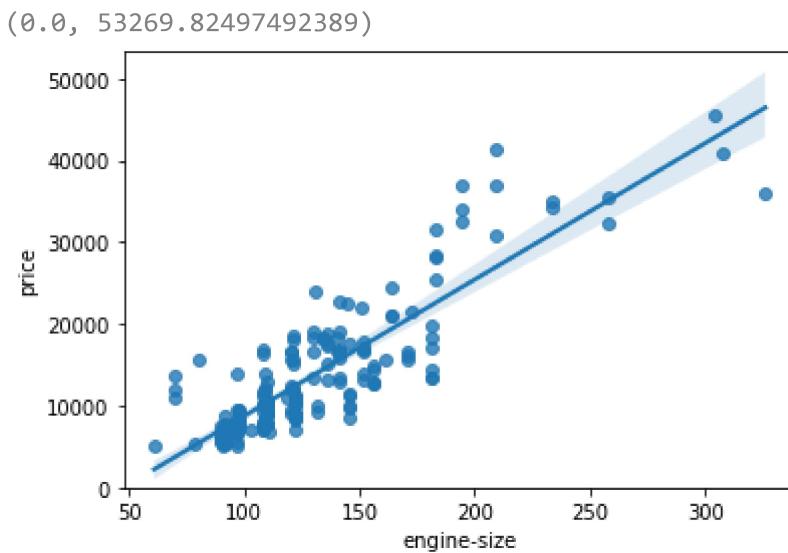
In order to start understanding the (linear) relationship between an individual variable and the price, we can use "regplot" which plots the scatterplot plus the fitted regression line for the data.

Let's see several examples of different linear relationships:

## Positive Linear Relationship

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

```
# 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 that it's approximately 0.87.

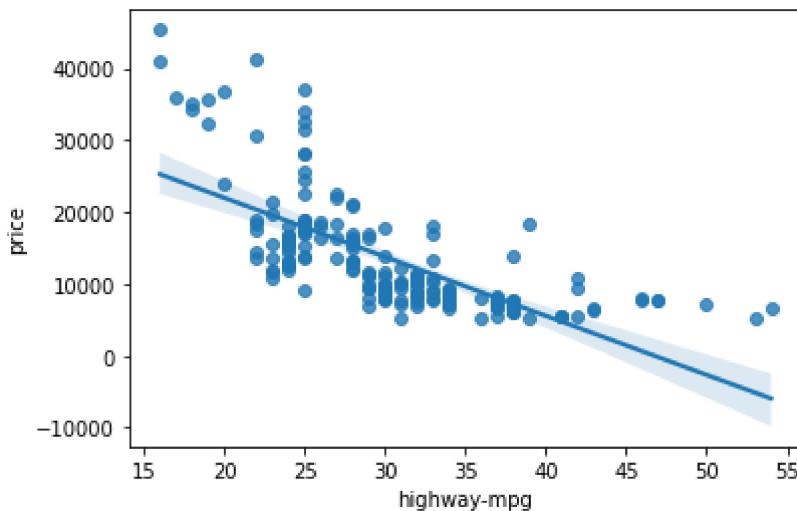
```
df[["engine-size", "price"]].corr()
```

engine-size      price      ⚡

Highway mpg is a potential predictor variable of price. Let's find the scatterplot of "highway-mpg" and "price".

```
sns.regplot(x="highway-mpg", y="price", data=df)
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7fb53882e350>
```



As 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.

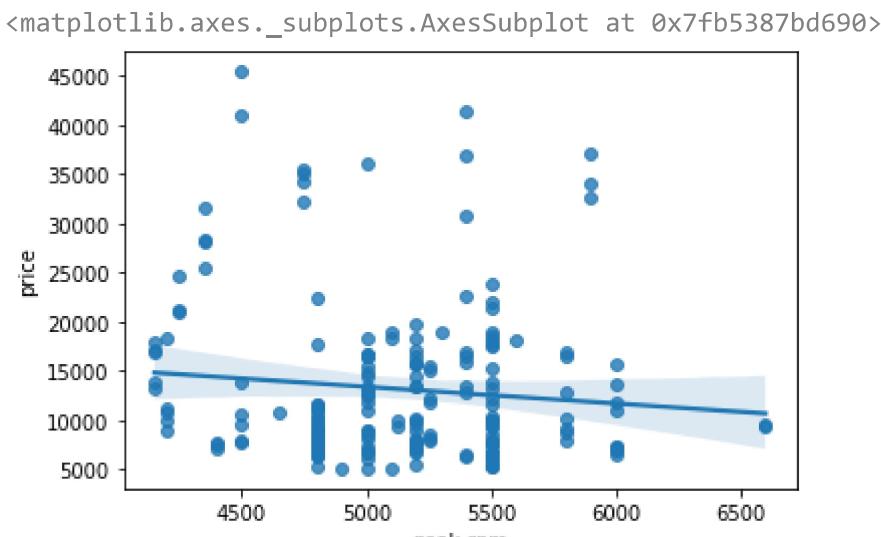
```
df[['highway-mpg', 'price']].corr()
```

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

## Weak Linear Relationship

Let's see if "peak-rpm" is a predictor variable of "price".

```
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.

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

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

## Question 3 a):

Find the correlation between x="stroke" and y="price".

Hint: if you would like to select those columns, use the following syntax: df[["stroke","price"]].

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

	stroke	price	edit
stroke	1.00000	0.08231	edit
price	0.08231	1.00000	edit

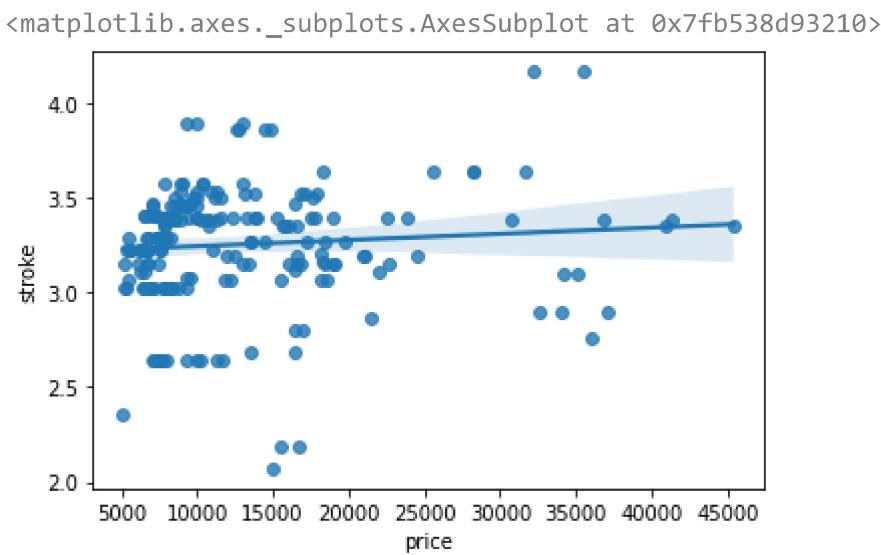
► Click here for the solution

## Question 3 b):

Given the correlation results between "price" and "stroke", do you expect a linear relationship?

Verify your results using the function "regplot()".

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



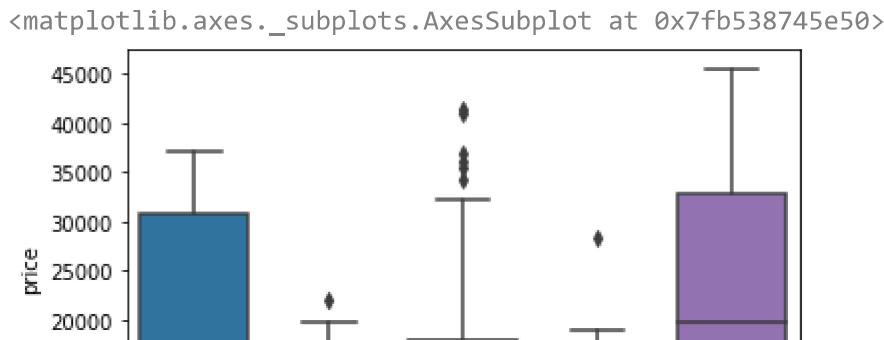
► Click here for the solution

## 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.

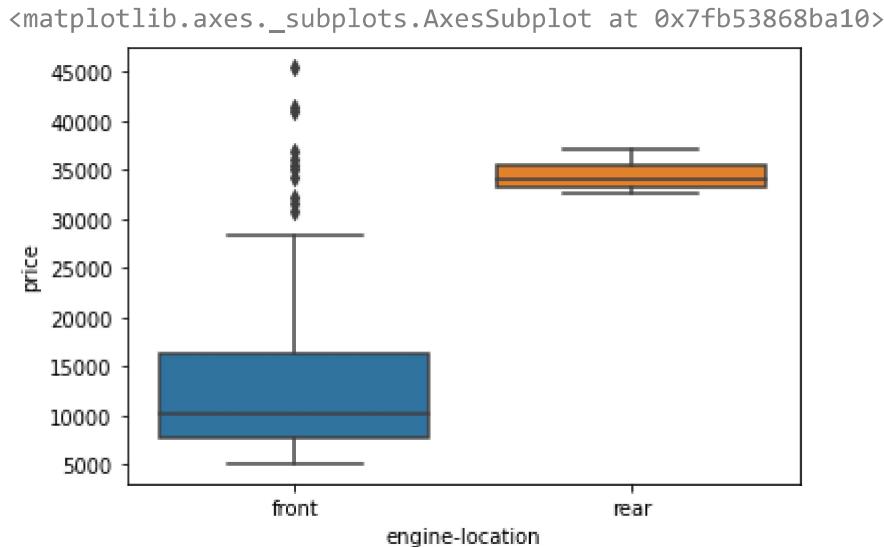
Let's look at the relationship between "body-style" and "price".

```
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, so body-style would not be a good predictor of price. Let's examine engine "engine-location" and "price":

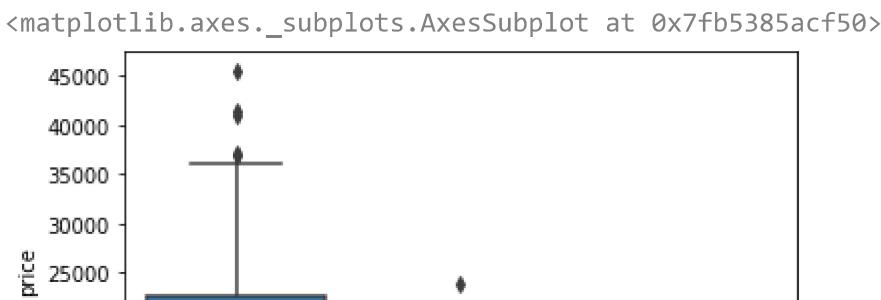
```
body-style
sns.boxplot(x="engine-location", y="price", data=df)
```



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".

```
# 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:

```
df.describe()
```

	symboling	normalized-losses	wheel-base	length	width	height	curb-weight
<b>count</b>	201.000000	201.000000	201.000000	201.000000	201.000000	201.000000	201.000000
<b>mean</b>	0.840796	122.000000	98.797015	0.837102	0.915126	53.766667	2555.666667
<b>std</b>	1.021000	31.000000	6.000000	0.060010	0.000107	0.117000	547.000000

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

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

	make	aspiration	num-of-doors	body-style	drive-wheels	engine-location	engine-type	num-of-cylinders	fuel-system
<b>count</b>	201	201	201	201	201	201	201	201	201
<b>unique</b>	22		2	2	5	3	2	6	7
<b>top</b>	toyota		std	four	sedan	fwd	front	ohc	four
<b>freq</b>	32		165	115	94	118	198	145	157

## 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']]`.

```
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:

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

drive-wheels	
fwd	118
rwd	75
-	-
-	-

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'.

```
drive_wheels_counts = df['drive-wheels'].value_counts().to_frame()
drive_wheels_counts.rename(columns={'drive-wheels': 'value_counts'}, inplace=True)
drive_wheels_counts
```

value_counts	
fwd	118
rwd	75
4wd	8

Now let's rename the index to 'drive-wheels':

```
drive_wheels_counts.index.name = 'drive-wheels'
drive_wheels_counts
```

value_counts	
drive-wheels	
fwd	118
rwd	75
4wd	8

We can repeat the above process for the variable 'engine-location'.

```
# engine-location as variable
engine_loc_counts = df['engine-location'].value_counts().to_frame()
engine_loc_counts.rename(columns={'engine-location': 'value_counts'}, inplace=True)
engine_loc_counts.index.name = 'engine-location'
engine_loc_counts.head(10)
```

value\_counts 

### engine-location

After examining the value counts of the engine location, we see that engine location would not be a good predictor variable for the price. This is because we only have three cars with a rear engine and 198 with an engine in the front, so this result is skewed. Thus, we are not able to draw any conclusions about the engine location.

## 4. Basics of Grouping

The "groupby" method groups data by different categories. The data is grouped based on one or several variables, and analysis is performed on the individual groups.

For example, let's group by the variable "drive-wheels". We see that there are 3 different categories of drive wheels.

```
df['drive-wheels'].unique()

array(['rwd', 'fwd', '4wd'], dtype=object)
```

If we want to know, on average, which type of drive wheel is most valuable, we can group "drive-wheels" and then average them.

We can select the columns 'drive-wheels', 'body-style' and 'price', then assign it to the variable "df\_group\_one".

```
df_group_one = df[['drive-wheels','body-style','price']]
```

We can then calculate the average price for each of the different categories of data.

```
# grouping results
df_group_one = df_group_one.groupby(['drive-wheels'],as_index=False).mean()
df_group_one
```

	drive-wheels	price	
0	4wd	10241.000000	
1	fwd	9244.779661	
2	rwd	19757.613333	

From our data, it seems rear-wheel drive vehicles are, on average, the most expensive, while 4-wheel and front-wheel are approximately the same in price.

You can also group by multiple variables. For example, let's group by both 'drive-wheels' and 'body-style'. This groups the dataframe by the unique combination of 'drive-wheels' and 'body-style'. We can store the results in the variable 'grouped\_test1'.

```
# grouping results
```

```
df_gptest = df[['drive-wheels','body-style','price']]
```

```
grouped_test1 = df_gptest.groupby(['drive-wheels','body-style'],as_index=False).mean()
grouped_test1
```

	drive-wheels	body-style	price	edit
0	4wd	hatchback	7603.000000	
1	4wd	sedan	12647.333333	
2	4wd	wagon	9095.750000	
3	fwd	convertible	11595.000000	
4	fwd	hardtop	8249.000000	
5	fwd	hatchback	8396.387755	
6	fwd	sedan	9811.800000	
7	fwd	wagon	9997.333333	
8	rwd	convertible	23949.600000	
9	rwd	hardtop	24202.714286	
10	rwd	hatchback	14337.777778	
11	rwd	sedan	21711.833333	
12	rwd	wagon	16994.222222	

This grouped data is much easier to visualize when it is made into a pivot table. A pivot table is like an Excel spreadsheet, with one variable along the column and another along the row. We can convert the dataframe to a pivot table using the method "pivot" to create a pivot table from the groups.

In this case, we will leave the drive-wheels variable as the rows of the table, and pivot body-style to become the columns of the table:

```
grouped_pivot = grouped_test1.pivot(index='drive-wheels',columns='body-style')
grouped_pivot
```



price					
body-style	convertible	hardtop	hatchback	sedan	wagon
drive-wheels					
4wd	NaN	NaN	7603.000000	12647.333333	9095.750000
fwd	11595.0	8249.000000	8396.387755	9811.800000	9997.333333
rwd	23949.6	24202.714286	14337.777778	21711.833333	16994.222222

Often, we won't have data for some of the pivot cells. We can fill these missing cells with the value 0, but any other value could potentially be used as well. It should be mentioned that missing data is quite a complex subject and is an entire course on its own.

```
grouped_pivot = grouped_pivot.fillna(0) #fill missing values with 0
grouped_pivot
```



price					
body-style	convertible	hardtop	hatchback	sedan	wagon
drive-wheels					
4wd	0.0	0.000000	7603.000000	12647.333333	9095.750000
fwd	11595.0	8249.000000	8396.387755	9811.800000	9997.333333
rwd	23949.6	24202.714286	14337.777778	21711.833333	16994.222222

## Question 4:

Use the "groupby" function to find the average "price" of each car based on "body-style".

```
# Write your code below and press Shift+Enter to execute
df_group_one = df[['drive-wheels','body-style','price']]
df_group_one = df_group_one.groupby(['body-style'],as_index=False).mean()
df_group_one
```

body-style	price	
0 convertible	21890.500000	edit

► Click here for the solution

2 hatchback 9957 111176

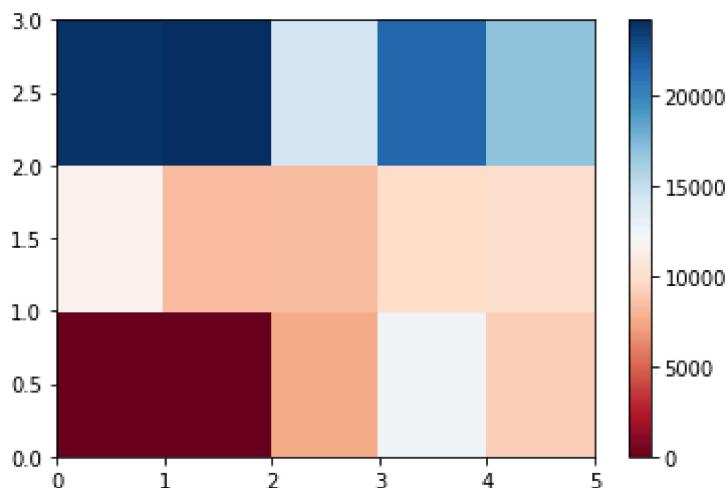
If you did not import "pyplot", let's do it again.

```
import matplotlib.pyplot as plt
%matplotlib inline
```

## Variables: Drive Wheels and Body Style vs. Price

Let's use a heat map to visualize the relationship between Body Style vs Price.

```
#use the grouped results
plt.pcolor(grouped_pivot, cmap='RdBu')
plt.colorbar()
plt.show()
```



The heatmap plots the target variable (price) proportional to colour with respect to the variables 'drive-wheel' and 'body-style' on the vertical and horizontal axis, respectively. This allows us to visualize how the price is related to 'drive-wheel' and 'body-style'.

The default labels convey no useful information to us. Let's change that:

```
fig, ax = plt.subplots()
im = ax.pcolor(grouped_pivot, cmap='RdBu')

#label names
row_labels = grouped_pivot.columns.levels[1]
https://colab.research.google.com/drive/1EKOPEcHgjzXKmxZPlmRYLbe5qjsUBm2t#scrollTo=L9iqAeoer3EP&printMode=true
```

```

col_labels = grouped_pivot.index

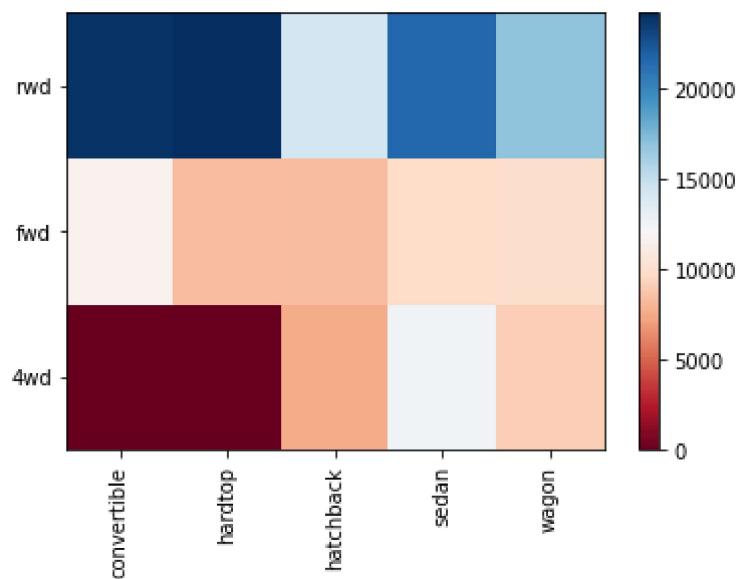
#move ticks and labels to the center
ax.set_xticks(np.arange(grouped_pivot.shape[1]) + 0.5, minor=False)
ax.set_yticks(np.arange(grouped_pivot.shape[0]) + 0.5, minor=False)

#insert labels
ax.set_xticklabels(row_labels, minor=False)
ax.set_yticklabels(col_labels, minor=False)

#rotate label if too long
plt.xticks(rotation=90)

fig.colorbar(im)
plt.show()

```



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. Correlation does not imply causation. Determining correlation is much simpler than determining causation as causation may require independent experimentation.

## Pearson Correlation

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

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

```
df.corr()
```

	<b>symboling</b>	<b>normalized-losses</b>	<b>wheel-base</b>	<b>length</b>	<b>width</b>	<b>height</b>	<b>curb-weight</b>
<b>symboling</b>	1.000000	0.466264	-0.535987	-0.365404	-0.242423	-0.550160	-0.233118
<b>normalized-losses</b>	0.466264	1.000000	-0.056661	0.019424	0.086802	-0.373737	0.099404
<b>wheel-base</b>	-0.535987	-0.056661	1.000000	0.876024	0.814507	0.590742	0.782097
<b>length</b>	-0.365404	0.019424	0.876024	1.000000	0.857170	0.492063	0.880665
<b>width</b>	-0.242423	0.086802	0.814507	0.857170	1.000000	0.306002	0.866201

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.

```
price      -0.082391      0.133999      0.584642      0.690628      0.751265      0.135486      0.834415
```

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

```
from scipy import stats
```

## Wheel-Base vs. Price

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

```
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 =", p_value)
```

The Pearson Correlation Coefficient is 0.584641822265508 with a P-value of P = 8.07648

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 ( $\sim 0.585$ ).

## Horsepower vs. Price

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

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

The Pearson Correlation Coefficient is 0.8095745670036559 with a P-value of P = 6.369



### Conclusion:

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

## Length vs. Price

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

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

The Pearson Correlation Coefficient is 0.6906283804483638 with a P-value of P = 8.016



### Conclusion:

Since the p-value is  $< 0.001$ , the correlation between length and price is statistically significant, and the linear relationship is moderately strong ( $\sim 0.691$ ).

## Width vs. Price

Let's calculate the Pearson Correlation Coefficient and P-value of 'width' and 'price':

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

```
The Pearson Correlation Coefficient is 0.7512653440522673 with a P-value of P = 9.2003
```

## Conclusion:

Since the p-value is < 0.001, the correlation between width and price is statistically significant, and the linear relationship is quite strong (~0.751).

### ▼ Curb-Weight vs. Price

Let's calculate the Pearson Correlation Coefficient and P-value of 'curb-weight' and 'price':

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

```
The Pearson Correlation Coefficient is 0.8344145257702843 with a P-value of P = 2.189
```

## Conclusion:

Since the p-value is < 0.001, the correlation between curb-weight and price is statistically significant, and the linear relationship is quite strong (~0.834).

## Engine-Size vs. Price

Let's calculate the Pearson Correlation Coefficient and P-value of 'engine-size' and 'price':

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

```
The Pearson Correlation Coefficient is 0.8723351674455185 with a P-value of P = 9.2654
```

## Conclusion:

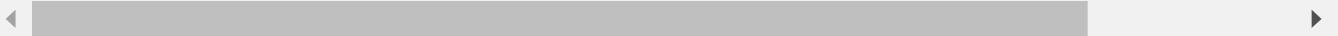
Since the p-value is < 0.001, the correlation between engine-size and price is statistically significant, and the linear relationship is very strong (~0.872).

## Bore vs. Price

Let's calculate the Pearson Correlation Coefficient and P-value of 'bore' and 'price':

```
pearson_coef, p_value = stats.pearsonr(df['bore'], df['price'])
print("The Pearson Correlation Coefficient is", pearson_coef, " with a P-value of P = ", p_\n
```

The Pearson Correlation Coefficient is 0.5431553832626602 with a P-value of P = 8.04



Conclusion:

Since the p-value is  $< 0.001$ , the correlation between bore and price is statistically significant, but the linear relationship is only moderate ( $\sim 0.521$ ).

We can relate the process for each 'city-mpg' and 'highway-mpg':

## City-mpg vs. Price

```
pearson_coef, p_value = stats.pearsonr(df['city-mpg'], df['price'])
print("The Pearson Correlation Coefficient is", pearson_coef, " with a P-value of P = ", p_\n
```

The Pearson Correlation Coefficient is -0.6865710067844678 with a P-value of P = 2.32



Conclusion:

Since the p-value is  $< 0.001$ , the correlation between city-mpg and price is statistically significant, and the coefficient of about -0.687 shows that the relationship is negative and moderately strong.

## Highway-mpg vs. Price

```
pearson_coef, p_value = stats.pearsonr(df['highway-mpg'], df['price'])
print("The Pearson Correlation Coefficient is", pearson_coef, " with a P-value of P = ", p_\n
```

The Pearson Correlation Coefficient is -0.704692265058953 with a P-value of P = 1.749



▼ Conclusion:

Since the p-value is  $< 0.001$ , the correlation between highway-mpg and price is statistically significant, and the coefficient of about -0.705 shows that the relationship is negative and

moderately strong.

## 6. ANOVA

### ANOVA: Analysis of Variance

The Analysis of Variance (ANOVA) is a statistical method used to test whether there are significant differences between the means of two or more groups. ANOVA returns two parameters:

**F-test score:** ANOVA assumes the means of all groups are the same, calculates how much the actual means deviate from the assumption, and reports it as the F-test score. A larger score means there is a larger difference between the means.

**P-value:** P-value tells how statistically significant our calculated score value is.

If our price variable is strongly correlated with the variable we are analyzing, we expect ANOVA to return a sizeable F-test score and a small p-value.

### Drive Wheels

Since ANOVA analyzes the difference between different groups of the same variable, the groupby function will come in handy. Because the ANOVA algorithm averages the data automatically, we do not need to take the average before hand.

To see if different types of 'drive-wheels' impact 'price', we group the data.

```
grouped_test2=df_gptest[['drive-wheels', 'price']].groupby(['drive-wheels'])
grouped_test2.head(2)
```

drive-wheels	price	
0	rwd	13495.0
1	rwd	16500.0
3	fwd	13950.0
4	4wd	17450.0
5	fwd	15250.0
136	4wd	7603.0

	drive-wheels	body-style	price	edit
0	rwd	convertible	13495.0	
1	rwd	convertible	16500.0	
2	rwd	hatchback	16500.0	
3	fwd	sedan	13950.0	
4	4wd	sedan	17450.0	
...	...	...	...	
196	rwd	sedan	16845.0	
197	rwd	sedan	19045.0	
198	rwd	sedan	21485.0	
199	rwd	sedan	22470.0	
200	rwd	sedan	22625.0	

201 rows × 3 columns

We can obtain the values of the method group using the method "get\_group".

```
grouped_test2.get_group('4wd')['price']
```

```
4      17450.0
136     7603.0
140     9233.0
141    11259.0
144     8013.0
145    11694.0
150     7898.0
151     8778.0
Name: price, dtype: float64
```

We can use the function 'f\_oneway' in the module 'stats' to obtain the **F-test score** and **P-value**.

```
# ANOVA
f_val, p_val = stats.f_oneway(grouped_test2.get_group('fwd')['price'], grouped_test2.get_group('rwd')['price'])

print("ANOVA results: F=", f_val, ", P =", p_val)

ANOVA results: F= 67.95406500780399 , P = 3.3945443577151245e-23
```

This is a great result with a large F-test score showing a strong correlation and a P-value of almost 0 implying almost certain statistical significance. But does this mean all three tested groups are all this highly correlated?

I can't examine them separately.

## ▼ fwd and rwd

```
f_val, p_val = stats.f_oneway(grouped_test2.get_group('fwd')['price'], grouped_test2.get_group('rwd')['price'])

print( "ANOVA results: F=", f_val, ", P =", p_val )

ANOVA results: F= 130.5533160959111 , P = 2.2355306355677845e-23
```

Let's examine the other groups.

## ▼ 4wd and rwd

```
f_val, p_val = stats.f_oneway(grouped_test2.get_group('4wd')['price'], grouped_test2.get_group('rwd')['price'])

print( "ANOVA results: F=", f_val, ", P =", p_val)

ANOVA results: F= 8.580681368924756 , P = 0.004411492211225333
```

4wd and fwd

```
f_val, p_val = stats.f_oneway(grouped_test2.get_group('4wd')['price'], grouped_test2.get_group('fwd')['price'])

print("ANOVA results: F=", f_val, ", P =", p_val)

ANOVA results: F= 0.665465750252303 , P = 0.41620116697845666
```

## Conclusion: Important Variables

We now have a better idea of what our data looks like and which variables are important to take into account when predicting the car price. We have narrowed it down to the following variables:

Continuous numerical variables:

- Length
- Width
- Curb-weight

- Engine-size
- Horsepower
- City-mpg
- Highway-mpg
- Wheel-base
- Bore

Categorical variables:

- Drive-wheels

As we now move into building machine learning models to automate our analysis, feeding the model with variables that meaningfully affect our target variable will improve our model's prediction performance.