**Skills Network**

# Multiple Linear Regression

Estimated time needed: **15** minutes

## Objectives

After completing this lab you will be able to:

- Use scikit-learn to implement Multiple Linear Regression
- Create a model, train it, test it and use the model

# Table of contents

---

## Importing Needed packages

```
In [6]:    import matplotlib.pyplot as plt
           import pandas as pd
           import pylab as pl
           import numpy as np
           %matplotlib inline
```

## Downloading Data

To download the data, we will use !wget to download it from IBM Object Storage.

```
In [7]:    !wget -O FuelConsumption.csv https://cf-courses-data.s3.us.cloud-object-storage.appdoma
```

```
--2024-09-16 15:43:52--  https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cl
oud/IBMDeveloperSkillsNetwork-ML0101EN-SkillsNetwork/labs/Module%202/data/FuelConsumptio
```

```
nCo2.csv
Resolving cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud (cf-courses-data.s
3.us.cloud-object-storage.appdomain.cloud)... 169.63.118.104, 169.63.118.104
Connecting to cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud (cf-courses-dat
a.s3.us.cloud-object-storage.appdomain.cloud)|169.63.118.104|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 72629 (71K) [text/csv]
Saving to: 'FuelConsumption.csv'

FuelConsumption.csv 100%[===================>]  70.93K  --.-KB/s    in 0.002s

2024-09-16 15:43:52 (45.2 MB/s) - 'FuelConsumption.csv' saved [72629/72629]
```

# Understanding the Data

## FuelConsumption.csv:

We have downloaded a fuel consumption dataset, **FuelConsumption.csv** , which contains model-specific fuel consumption ratings and estimated carbon dioxide emissions for new light-duty vehicles for retail sale in Canada. Dataset source

- **MODELYEAR** e.g. 2014
- **MAKE** e.g. Acura
- **MODEL** e.g. ILX
- **VEHICLE CLASS** e.g. SUV
- **ENGINE SIZE** e.g. 4.7
- **CYLINDERS** e.g 6
- **TRANSMISSION** e.g. A6
- **FUELTYPE** e.g. z
- **FUEL CONSUMPTION in CITY(L/100 km)** e.g. 9.9
- **FUEL CONSUMPTION in HWY (L/100 km)** e.g. 8.9
- **FUEL CONSUMPTION COMB (L/100 km)** e.g. 9.2
- **CO2 EMISSIONS (g/km)** e.g. 182 --> low --> 0

# Reading the data in

In [8]:
```python
df = pd.read_csv("FuelConsumption.csv")

# take a look at the dataset
df.head()
```

Out[8]:

| | MODELYEAR | MAKE | MODEL | VEHICLECLASS | ENGINESIZE | CYLINDERS | TRANSMISSION | FUELTYPE | |
|---|---|---|---|---|---|---|---|---|---|
| **0** | 2014 | ACURA | ILX | COMPACT | 2.0 | 4 | AS5 | Z | |
| **1** | 2014 | ACURA | ILX | COMPACT | 2.4 | 4 | M6 | Z | |
| **2** | 2014 | ACURA | ILX HYBRID | COMPACT | 1.5 | 4 | AV7 | Z | |

| | MODELYEAR | MAKE | MODEL | VEHICLECLASS | ENGINESIZE | CYLINDERS | TRANSMISSION | FUELTYPE | |
|---|---|---|---|---|---|---|---|---|---|
| **3** | 2014 | ACURA | MDX 4WD | SUV - SMALL | 3.5 | 6 | AS6 | Z | |
| **4** | 2014 | ACURA | RDX AWD | SUV - SMALL | 3.5 | 6 | AS6 | Z | |

Let's select some features that we want to use for regression.

In [9]:
```python
cdf = df[['ENGINESIZE','CYLINDERS','FUELCONSUMPTION_CITY','FUELCONSUMPTION_HWY','FUELCO
cdf.head(9)
```

Out[9]:

| | ENGINESIZE | CYLINDERS | FUELCONSUMPTION_CITY | FUELCONSUMPTION_HWY | FUELCONSUMPTION_C( |
|---|---|---|---|---|---|
| **0** | 2.0 | 4 | 9.9 | 6.7 | |
| **1** | 2.4 | 4 | 11.2 | 7.7 | |
| **2** | 1.5 | 4 | 6.0 | 5.8 | |
| **3** | 3.5 | 6 | 12.7 | 9.1 | |
| **4** | 3.5 | 6 | 12.1 | 8.7 | |
| **5** | 3.5 | 6 | 11.9 | 7.7 | |
| **6** | 3.5 | 6 | 11.8 | 8.1 | |
| **7** | 3.7 | 6 | 12.8 | 9.0 | |
| **8** | 3.7 | 6 | 13.4 | 9.5 | |

Let's plot Emission values with respect to Engine size:

In [10]:
```python
plt.scatter(cdf.ENGINESIZE, cdf.CO2EMISSIONS,  color='blue')
plt.xlabel("Engine size")
plt.ylabel("Emission")
plt.show()
```

## Creating train and test dataset

Train/Test Split involves splitting the dataset into training and testing sets respectively, which are mutually exclusive. After which, you train with the training set and test with the testing set. This will provide a more accurate evaluation on out-of-sample accuracy because the testing dataset is not part of the dataset that have been used to train the model. Therefore, it gives us a better understanding of how well our model generalizes on new data.
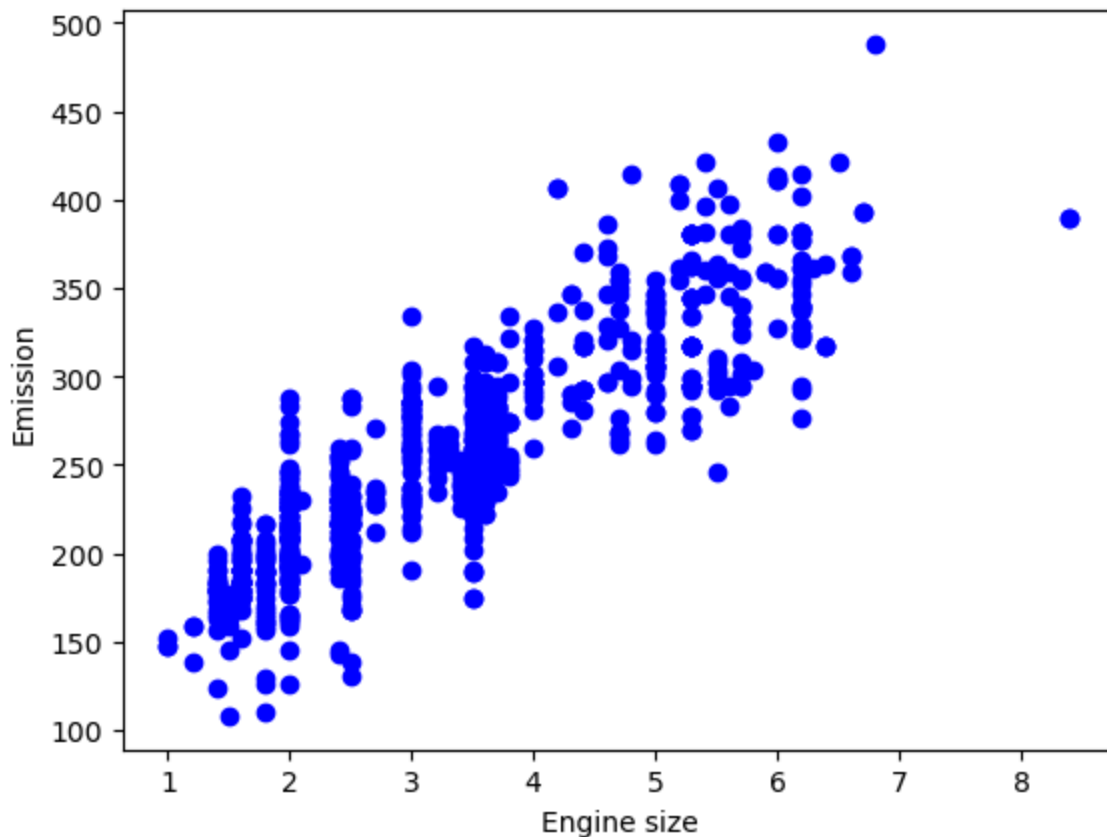
We know the outcome of each data point in the testing dataset, making it great to test with! Since this data has not been used to train the model, the model has no knowledge of the outcome of these data points. So, in essence, it is truly an out-of-sample testing.

Let's split our dataset into train and test sets. Around 80% of the entire dataset will be used for training and 20% for testing. We create a mask to select random rows using the **np.random.rand()** function:

In [11]:
```python
msk = np.random.rand(len(df)) < 0.8
train = cdf[msk]
test = cdf[~msk]
```

## Train data distribution

In [12]:
```python
plt.scatter(train.ENGINESIZE, train.CO2EMISSIONS,  color='blue')
plt.xlabel("Engine size")
plt.ylabel("Emission")
plt.show()
```

# Multiple Regression Model

In reality, there are multiple variables that impact the co2emission. When more than one independent variable is present, the process is called multiple linear regression. An example of multiple linear regression is predicting co2emission using the features FUELCONSUMPTION_COMB, EngineSize and Cylinders of cars. The good thing here is that multiple linear regression model is the extension of the simple linear regression model.

In [16]:
```python
from sklearn import linear_model
regr = linear_model.LinearRegression()
x = np.asanyarray(train[['ENGINESIZE','CYLINDERS','FUELCONSUMPTION_COMB']])
y = np.asanyarray(train[['CO2EMISSIONS']])
regr.fit (x, y)
# The coefficients
print ('Coefficients: ', regr.coef_)
```

```
Coefficients:  [[10.82104577  6.92142797  9.71244165]]
```

As mentioned before, **Coefficient** and **Intercept** are the parameters of the fitted line. Given that it is a multiple linear regression model with 3 parameters and that the parameters are the intercept and coefficients of the hyperplane, sklearn can estimate them from our data. Scikit-learn uses plain Ordinary Least Squares method to solve this problem.

### Ordinary Least Squares (OLS)

OLS is a method for estimating the unknown parameters in a linear regression model. OLS chooses the parameters of a linear function of a set of explanatory variables by minimizing the sum of the

squares of the differences between the target dependent variable and those predicted by the linear function. In other words, it tries to minimizes the sum of squared errors (SSE) or mean squared error (MSE) between the target variable (y) and our predicted output ($\hat{y}$) over all samples in the dataset.

OLS can find the best parameters using of the following methods:

- Solving the model parameters analytically using closed-form equations
- Using an optimization algorithm (Gradient Descent, Stochastic Gradient Descent, Newton's Method, etc.)

# Prediction

In [14]:
```python
y_hat= regr.predict(test[['ENGINESIZE','CYLINDERS','FUELCONSUMPTION_COMB']])
x = np.asanyarray(test[['ENGINESIZE','CYLINDERS','FUELCONSUMPTION_COMB']])
y = np.asanyarray(test[['CO2EMISSIONS']])
print("Mean Squared Error (MSE) : %.2f"
      % np.mean((y_hat - y) ** 2))

# Explained variance score: 1 is perfect prediction
print('Variance score: %.2f' % regr.score(x, y))
```

```
Mean Squared Error (MSE) : 642.64
Variance score: 0.86
```

**Explained variance regression score:**

Let $\hat{y}$ be the estimated target output, y the corresponding (correct) target output, and Var be the Variance (the square of the standard deviation). Then the explained variance is estimated as follows:

$$\text{explainedVariance}(y, \hat{y}) = 1 - \frac{Var\{y-\hat{y}\}}{Var\{y\}}$$

The best possible score is 1.0, the lower values are worse.

# Practice

Try to use a multiple linear regression with the same dataset, but this time use FUELCONSUMPTION_CITY and FUELCONSUMPTION_HWY instead of FUELCONSUMPTION_COMB. Does it result in better accuracy?

In [15]:
```python
regr = linear_model.LinearRegression()
x = np.asanyarray(train[['ENGINESIZE','CYLINDERS','FUELCONSUMPTION_CITY','FUELCONSUMPTI
y = np.asanyarray(train[['CO2EMISSIONS']])
regr.fit (x, y)
print ('Coefficients: ', regr.coef_)
y_= regr.predict(test[['ENGINESIZE','CYLINDERS','FUELCONSUMPTION_CITY','FUELCONSUMPTION
x = np.asanyarray(test[['ENGINESIZE','CYLINDERS','FUELCONSUMPTION_CITY','FUELCONSUMPTIO
y = np.asanyarray(test[['CO2EMISSIONS']])
print("Residual sum of squares: %.2f"% np.mean((y_ - y) ** 2))
print('Variance score: %.2f' % regr.score(x, y))
```

```
Coefficients:  [[10.84138186  6.83784697  5.55570721  4.09184974]]
Residual sum of squares: 642.43
Variance score: 0.86
```

▶ Click here for the solution

# Thank you for completing this lab!

# Author

Saeed Aghabozorgi

## Other Contributors

Joseph Santarcangelo

## © IBM Corporation 2020. All rights reserved.

<!--

# Change Log

| Date (YYYY-MM-DD) | Version | Changed By | Change Description |
|---|---|---|---|
| 2020-11-03 | 2.1 | Lakshmi | Made changes in URL |
| 2020-08-27 | 2.0 | Lavanya | Moved lab to course repo in GitLab |

--!>

```
In [ ]:
```