

ML0101EN-Reg-Multiple-Linear-Regression-Co2-py-v1

December 5, 2018

```
#  
Multiple Linear Regression
```

About this Notebook In this notebook, we learn how to use scikit-learn to implement Multiple linear regression. We download a dataset that is related to fuel consumption and Carbon dioxide emission of cars. Then, we split our data into training and test sets, create a model using training set, Evaluate your model using test set, and finally use model to predict unknown value

0.0.1 Importing Needed packages

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

0.0.2 Downloading Data

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

```
In [ ]: !wget -O FuelConsumption.csv https://s3-api.us-geo.objectstorage.softlayer.net/cf-course
```

Did you know? When it comes to Machine Learning, you will likely be working with large datasets. As a business, where can you host your data? IBM is offering a unique opportunity for businesses, with 10 Tb of IBM Cloud Object Storage: [Sign up now for free](#)

0.1 Understanding the Data

0.1.1 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

0.2 Reading the data in

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

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

Lets select some features that we want to use for regression.

```
In [ ]: cdf = df[['ENGINE SIZE', 'CYLINDERS', 'FUELCONSUMPTION_CITY', 'FUELCONSUMPTION_HWY', 'FUELCONSUMPTION_COMB', 'CO2EMISSIONS']]
cdf.head(9)
```

Lets plot Emission values with respect to Engine size:

```
In [ ]: plt.scatter(cdf.ENGINE SIZE, 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 data. It is more realistic for real world problems.

This means that we know the outcome of each data point in this dataset, making it great to test with! And 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's truly an out-of-sample testing.

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

Train data distribution

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

0.3 Multiple Regression Model

In reality, there are multiple variables that predict the Co2emission. When more than one independent variable is present, the process is called multiple linear regression. For example, predicting co2emission using FUELCONSUMPTION_COMB, EngineSize and Cylinders of cars. The good thing here is that Multiple linear regression is the extension of simple linear regression model.

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

As mentioned before, **Coefficient** and **Intercept**, are the parameters of the fit line. Given that it is a multiple linear regression, with 3 parameters, and knowing that the parameters are the intercept and coefficients of 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 minimize 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.)

0.3.1 Prediction

```
In [ ]: y_hat= regr.predict(test[['ENGINE_SIZE', 'CYLINDERS', 'FUELCONSUMPTION_COMB']])
        x = np.asanyarray(test[['ENGINE_SIZE', 'CYLINDERS', 'FUELCONSUMPTION_COMB']])
        y = np.asanyarray(test[['CO2EMISSIONS']])
        print("Residual sum of squares: %.2f"
              % np.mean((y_hat - y) ** 2))

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

explained variance regression score:

If \hat{y} is the estimated target output, y the corresponding (correct) target output, and Var is Variance, the square of the standard deviation, then the explained variance is estimated as follow:

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

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

0.4 Practice

Try to use a multiple linear regression with the same dataset but this time use **FUEL CONSUMPTION in CITY** and **FUEL CONSUMPTION in HWY** instead of **FUELCONSUMPTION_COMB**. Does it result in better accuracy?

```
In [ ]: # write your code here
```

Double-click **here** for the solution.

0.5 Want to learn more?

IBM SPSS Modeler is a comprehensive analytics platform that has many machine learning algorithms. It has been designed to bring predictive intelligence to decisions made by individuals, by groups, by systems – by your enterprise as a whole. A free trial is available through this course, available here: [SPSS Modeler](#).

Also, you can use Watson Studio to run these notebooks faster with bigger datasets. Watson Studio is IBM's leading cloud solution for data scientists, built by data scientists. With Jupyter notebooks, RStudio, Apache Spark and popular libraries pre-packaged in the cloud, Watson Studio enables data scientists to collaborate on their projects without having to install anything. Join the fast-growing community of Watson Studio users today with a free account at [Watson Studio](#)

0.5.1 Thanks for completing this lesson!

Notebook created by: Saeed Aghabozorgi

Copyright © 2018 [Cognitive Class](#). This notebook and its source code are released under the terms of the [MIT License](#).