

ML0101EN-Reg-Simple-Linear-Regression-Co2

May 26, 2022

1 Simple Linear Regression

Estimated time needed: **15** minutes

1.1 Objectives

After completing this lab you will be able to:

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

1.1.1 Importing Needed packages

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

1.1.2 Downloading Data

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

```
[ ]: !wget -O FuelConsumption.csv https://cf-courses-data.s3.us.cloud-object-storage.
↪appdomain.cloud/IBMDeveloperSkillsNetwork-ML0101EN-SkillsNetwork/labs/
↪Module%202/data/FuelConsumptionCo2.csv
```

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](#)

1.2 Understanding the Data

1.2.1 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
- **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

1.3 Reading the data in

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

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

1.3.1 Data Exploration

Let's first have a descriptive exploration on our data.

```
[ ]: # summarize the data
df.describe()
```

Let's select some features to explore more.

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

We can plot each of these features:

```
[ ]: viz = cdf[['CYLINDERS', 'ENGINE SIZE', 'CO2EMISSIONS', 'FUELCONSUMPTION_COMB']]
viz.hist()
plt.show()
```

Now, let's plot each of these features against the Emission, to see how linear their relationship is:

```
[ ]: plt.scatter(cdf.FUELCONSUMPTION_COMB, cdf.CO2EMISSIONS, color='blue')
plt.xlabel("FUELCONSUMPTION_COMB")
plt.ylabel("Emission")
plt.show()
```

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

1.4 Practice

Plot **CYLINDER** vs the Emission, to see how linear is their relationship is:

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

[Click here for the solution](#)

```
plt.scatter(cdf.CYLINDERS, cdf.CO2EMISSIONS, color='blue')
plt.xlabel("Cylinders")
plt.ylabel("Emission")
plt.show()
```

Creating train and test dataset Train/Test Split involves splitting the dataset into training and testing sets that 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.

This means that 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. 80% of the entire dataset will be used for training and 20% for testing. We create a mask to select random rows using **np.random.rand()** function:

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

1.4.1 Simple Regression Model

Linear Regression fits a linear model with coefficients $B = (B_1, \dots, B_n)$ to minimize the 'residual sum of squares' between the actual value y in the dataset, and the predicted value \hat{y} using linear approximation.

Train data distribution

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

Modeling Using sklearn package to model data.

```
[ ]: from sklearn import linear_model
     regr = linear_model.LinearRegression()
     train_x = np.asanyarray(train[['ENGINE SIZE']])
     train_y = np.asanyarray(train[['CO2EMISSIONS']])
     regr.fit(train_x, train_y)
```

```
# The coefficients
print ('Coefficients: ', regr.coef_)
print ('Intercept: ',regr.intercept_)
```

As mentioned before, **Coefficient** and **Intercept** in the simple linear regression, are the parameters of the fit line. Given that it is a simple linear regression, with only 2 parameters, and knowing that the parameters are the intercept and slope of the line, sklearn can estimate them directly from our data. Notice that all of the data must be available to traverse and calculate the parameters.

Plot outputs We can plot the fit line over the data:

```
[ ]: plt.scatter(train.ENGINESIZE, train.CO2EMISSIONS, color='blue')
plt.plot(train_x, regr.coef_[0][0]*train_x + regr.intercept_[0], '-r')
plt.xlabel("Engine size")
plt.ylabel("Emission")
```

Evaluation We compare the actual values and predicted values to calculate the accuracy of a regression model. Evaluation metrics provide a key role in the development of a model, as it provides insight to areas that require improvement.

There are different model evaluation metrics, lets use MSE here to calculate the accuracy of our model based on the test set:

- Mean Absolute Error: It is the mean of the absolute value of the errors. This is the easiest of the metrics to understand since it's just average error.
- Mean Squared Error (MSE): Mean Squared Error (MSE) is the mean of the squared error. It's more popular than Mean Absolute Error because the focus is geared more towards large errors. This is due to the squared term exponentially increasing larger errors in comparison to smaller ones.
- Root Mean Squared Error (RMSE).
- R-squared is not an error, but rather a popular metric to measure the performance of your regression model. It represents how close the data points are to the fitted regression line. The higher the R-squared value, the better the model fits your data. The best possible score is 1.0 and it can be negative (because the model can be arbitrarily worse).

```
[ ]: from sklearn.metrics import r2_score

test_x = np.asanyarray(test[['ENGINE SIZE']])
test_y = np.asanyarray(test[['CO2 EMISSIONS']])
test_y_ = regr.predict(test_x)

print("Mean absolute error: %.2f" % np.mean(np.absolute(test_y_ - test_y)))
print("Residual sum of squares (MSE): %.2f" % np.mean((test_y_ - test_y) ** 2))
print("R2-score: %.2f" % r2_score(test_y , test_y_ ) )
```

1.5 Exercise

Lets see what the evaluation metrics are if we trained a regression model using the FUELCONSUMPTION_COMB feature.

Start by selecting FUELCONSUMPTION_COMB as the train_x data from the train dataframe, then select FUELCONSUMPTION_COMB as the test_x data from the test dataframe

```
[ ]: train_x = #ADD CODE

test_x = #ADD CODE
```

[Click here for the solution](#)

```
train_x = train[["FUELCONSUMPTION_COMB"]]

test_x = test[["FUELCONSUMPTION_COMB"]]
```

Now train a Logistic Regression Model using the train_x you created and the train_y created previously

```
[ ]: regr = linear_model.LinearRegression()

#ADD CODE
```

[Click here for the solution](#)

```
regr = linear_model.LinearRegression()

regr.fit(train_x, train_y)
```

Find the predictions using the model's predict function and the test_x data

```
[ ]: predictions = #ADD CODE
```

[Click here for the solution](#)

```
predictions = regr.predict(test_x)
```

Finally use the predictions and the test_y data and find the Mean Absolute Error value using the np.absolute and np.mean function like done previously

```
[ ]: #ADD CODE
```

[Click here for the solution](#)

```
print("Mean Absolute Error: %.2f" % np.mean(np.absolute(predictions - test_y)))
```

We can see that the MAE is much worse when we train using ENGINESIZE than FUELCONSUMPTION_COMB

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

1.5.1 Thank you for completing this lab!

1.6 Author

Saeed Aghabozorgi

1.6.1 Other Contributors

Joseph Santarcangelo

Azim Hirjani

1.7 Change Log

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

##

© IBM Corporation 2020. All rights reserved.