

Skills Network Labs

File Edit View Run Kernel Git Tabs Settings Help

ML0101EN-Reg-Polynomial

Python

IBM Developer SKILLS NETWORK

Polynomial Regression

Estimated time needed: 15 minutes

Objectives

After completing this lab you will be able to:

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

Table of contents

1. Downloading Data

2. Polynomial regression

3. Evaluation

4. Practice

Importing Needed packages

```
[ ]: 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 `lwget` to download it from IBM Object Storage.

```
[ ]: lwget -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](#)

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
- **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

```
[ ]: 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.

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

Lets plot Emission values with respect to Engine size:

```
[ ]: 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.

Support/Feedback

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

Polynomial regression

Did you know? IBM Watson Studio lets you build and deploy an AI solution, using the best of open source and IBM software and giving your team a single environment to work in. [Learn more here.](#)

Sometimes, the trend of data is not really linear, and looks curvy. In this case we can use Polynomial regression methods. In fact, many different regressions exist that can be used to fit whatever the dataset looks like, such as quadratic, cubic, and so on, and it can go on and on to infinite degrees.

In essence, we can call all of these, polynomial regression, where the relationship between the independent variable x and the dependent variable y is modeled as an n th degree polynomial in x . Lets say you want to have a polynomial regression (let's make 2 degree polynomial):

$$y = b + \theta_1 x + \theta_2 x^2$$

Now, the question is: how we can fit our data on this equation while we have only x values, such as **Engine Size**? Well, we can create a few additional features: 1, x , and x^2 .

PolynomialFeatures() function in Scikit-learn library, drives a new feature sets from the original feature set. That is, a matrix will be generated consisting of all polynomial combinations of the features with degree less than or equal to the specified degree. For example, lets say the original feature set has only one feature, *ENGINE SIZE*. Now, if we select the degree of the polynomial to be 2, then it generates 3 features, degree=0, degree=1 and degree=2:

```
[ ]: from sklearn.preprocessing import PolynomialFeatures
      from sklearn import linear_model
      train_x = np.asarray(train[['ENGINE SIZE']])
      train_y = np.asarray(train[['CO2EMISSIONS']])

      test_x = np.asarray(test[['ENGINE SIZE']])
      test_y = np.asarray(test[['CO2EMISSIONS']])

      poly = PolynomialFeatures(degree=2)
      train_x_poly = poly.fit_transform(train_x)
      train_x_poly
```

fit_transform takes our x values, and output a list of our data raised from power of 0 to power of 2 (since we set the degree of our polynomial to 2).

The equation and the sample example is displayed below.

$$\begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix} \rightarrow \begin{bmatrix} 1 & v_1 & v_1^2 \\ 1 & v_2 & v_2^2 \\ \vdots & \vdots & \vdots \\ 1 & v_n & v_n^2 \end{bmatrix}$$

$$\begin{bmatrix} 2. \\ 2.4 \\ 1.5 \\ \vdots \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 2. & 4. \\ 1 & 2.4 & 5.76 \\ 1 & 1.5 & 2.25 \\ \vdots & \vdots & \vdots \end{bmatrix}$$

It looks like feature sets for multiple linear regression analysis, right? Yes. It Does. Indeed, Polynomial regression is a special case of linear regression, with the main idea of how do you select your features. Just consider replacing the x with x_1 , x_1^2 with x_2 , and so on. Then the degree 2 equation would be turn into:

$$y = b + \theta_1 x_1 + \theta_2 x_2$$

Now, we can deal with it as 'linear regression' problem. Therefore, this polynomial regression is considered to be a special case of traditional multiple linear regression. So, you can use the same mechanism as linear regression to solve such a problems.

so we can use **LinearRegression()** function to solve it:

```
[ ]: clf = linear_model.LinearRegression()
      train_y_ = clf.fit(train_x_poly, train_y)
      # The coefficients
      print('Coefficients: ', clf.coef_)
      print('Intercept: ', clf.intercept_)
```

As mentioned before, **Coefficient** and **Intercept**, are the parameters of the fit curvy line. Given that it is a typical multiple linear regression, with 3 parameters, and knowing that the parameters are the intercept and coefficients of hyperplane, sklearn has estimated them from our new set of feature sets. Lets plot it:

```
[ ]: plt.scatter(train.ENGINE SIZE, train.CO2EMISSIONS, color='blue')
      XX = np.arange(0.0, 10.0, 0.1)
      yy = clf.intercept_[0] + clf.coef_[0][1]*XX + clf.coef_[0][2]*np.power(XX, 2)
      plt.plot(XX, yy, '-r')
      plt.xlabel("Engine size")
      plt.ylabel("Emission")
```

Evaluation

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

      test_x_poly = poly.fit_transform(test_x)
      test_y_ = clf.predict(test_x_poly)

      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_))
```

Practice

Try to use a polynomial regression with the dataset but this time with degree three (cubic). Does it result in better accuracy?

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

▼ Click here for the solution

```
poly3 = PolynomialFeatures(degree=3)
train_x_poly3 = poly3.fit_transform(train_x)
clf3 = linear_model.LinearRegression()
train_y3_ = clf3.fit(train_x_poly3, train_y)

# The coefficients
print('Coefficients: ', clf3.coef_)
```

```
print ('Intercept: ',clf3.intercept_)
plt.scatter(train.ENGINESIZE, train.CO2EMISSIONS, color='blue')
XX = np.arange(0.0, 10.0, 0.1)
yy = clf3.intercept_[0]+ clf3.coef_[0][1]*XX + clf3.coef_[0][2]*np.power(XX, 2) + clf3.coef_[0][3]*np.power(XX, 3)
plt.plot(XX, yy, '-r' )
plt.xlabel("Engine size")
plt.ylabel("Emission")
test_x_poly3 = poly3.fit_transform(test_x)
test_y3_ = clf3.predict(test_x_poly3)
print("Mean absolute error: %.2f" % np.mean(np.absolute(test_y3_ - test_y)))
print("Residual sum of squares (MSE): %.2f" % np.mean((test_y3_ - test_y) ** 2))
print("R2-score: %.2f" % r2_score(test_y3_ , test_y) )
```

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](#)

Thank you for completing this lab!

Author

Saeed Aghabozorgi

Other Contributors

[Joseph Santarcangelo](#)

Change Log

Date (YYYY-MM-DD)	Version	Changed By	Change Description
2020-11-04	2.2	Lakshmi	Made changes in markdown of equations
2020-11-03	2.1	Lakshmi	Made changes in URL
2020-08-27	2.0	Lavanya	Moved lab to course repo in GitLab

© IBM Corporation 2020. All rights reserved.