

REGRESSION

Simple Linear Regression

```
In [1]: import matplotlib.pyplot as plt
import pandas as pd
import pylab as pl
import numpy as np
%matplotlib inline
```

```
In [6]: df = pd.read_csv('https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBMDeveloperSkillsNetwork-
df.head(10)
```

	MODEL	YEAR	MAKE	MODEL	VEHICLECLASS	ENGINEIZE	CYLINDERS	TRANSMISSION	FUELTYPE	FUELCONSUMPTION_CITY	FUELCONSUMPTION_Hwy	CO2EMISSIONS
0	2014	ACURA	ILX	COMPACT	2.0	4	A55	Z		9.9		
1	2014	ACURA	ILX	COMPACT	2.4	4	M6	Z		11.2		
2	2014	ACURA	ILX	HYBRID	COMPACT	1.5	4	AV7	Z	6.0		
3	2014	ACURA	MDX	SUV - SMALL	3.5	6	A56	Z		12.7		
4	2014	ACURA	RDX	SUV - SMALL	3.5	6	A56	Z		12.1		

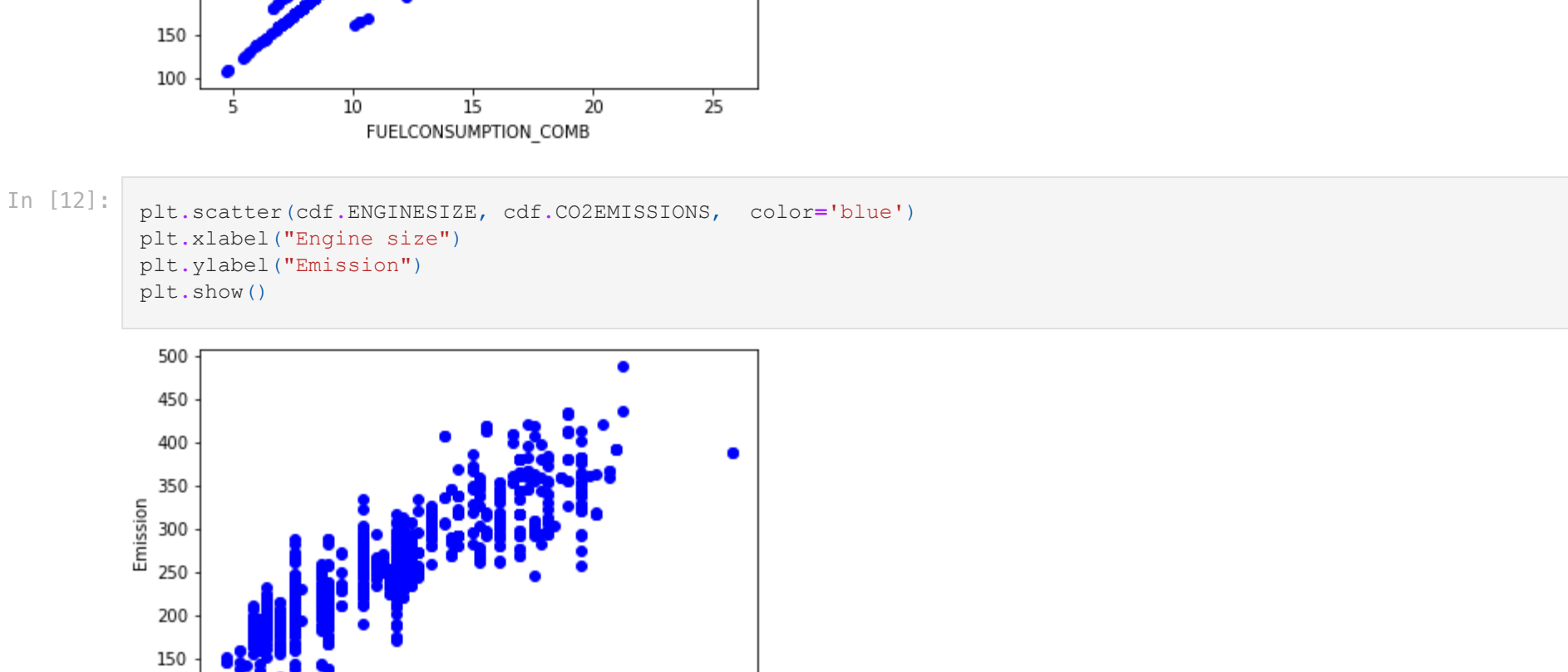
```
In [8]: df.describe()
```

	MODEL	YEAR	ENGINEIZE	CYLINDERS	FUELCONSUMPTION_CITY	FUELCONSUMPTION_Hwy	FUELCONSUMPTION_COMB	CO2EMISSIONS
count	10	10	10	10	10	10	10	10
mean	2014.0	3.346298	5.794752	4.101253	13.296532	9.474602	11.860881	3.485595
std	0.0	1.415895	1.797447	4.101253	2.794510	3.485595	4.700000	9.000000
min	2014.0	1.000000	3.000000	4.000000	4.600000	4.900000	4.700000	9.000000
25%	2014.0	2.000000	4.000000	10.250000	7.500000	9.000000	10.900000	9.000000
50%	2014.0	3.400000	6.000000	12.600000	8.800000	10.900000	13.350000	9.000000
75%	2014.0	4.300000	8.000000	15.500000	10.850000	13.350000	15.500000	13.350000
max	2014.0	8.400000	12.000000	30.200000	20.500000	25.800000	25.800000	25.800000

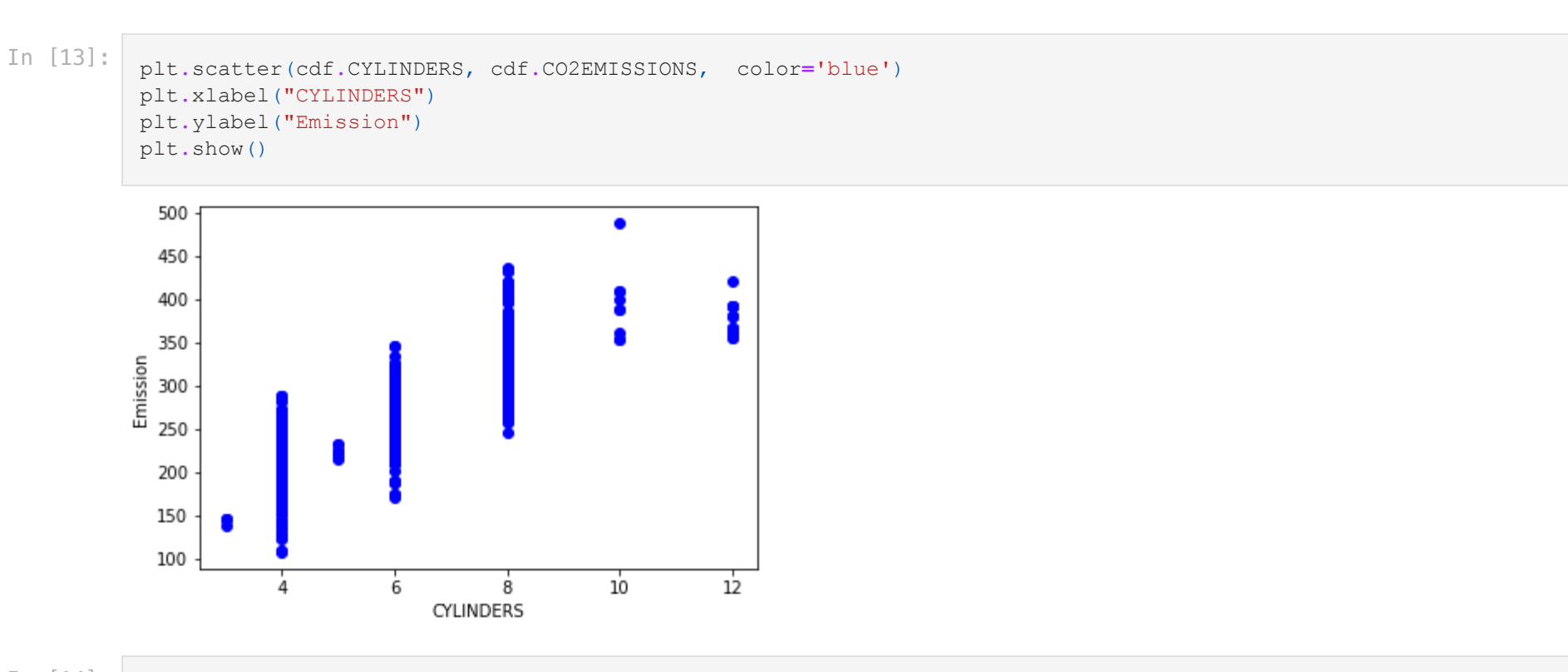
```
In [9]: #Exploring some of the features
cdf = df[['ENGINEIZE','CYLINDERS','FUELCONSUMPTION_COMB','CO2EMISSIONS']]
cdf.head(10)
```

	ENGINEIZE	CYLINDERS	FUELCONSUMPTION_COMB	CO2EMISSIONS
0	2.0	4	8.5	196
1	2.4	4	9.6	221
2	1.5	4	5.9	136
3	3.5	6	11.1	255
4	3.5	6	10.6	244
5	3.5	6	10.1	230
6	3.5	6	10.1	232
7	3.7	6	11.1	255
8	3.7	6	11.6	267
9	2.4	4	9.2	212

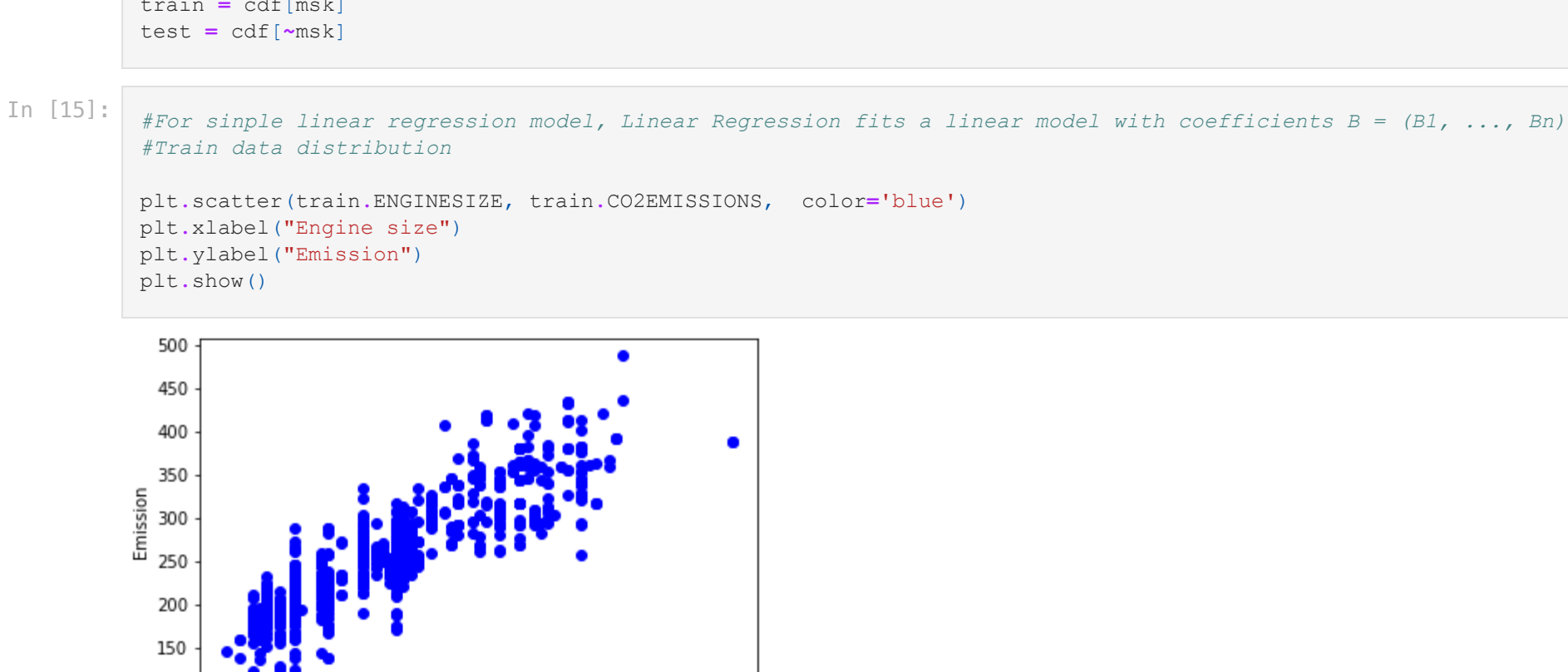
```
In [10]: # visualizing the above features by plotting them
vis = cdf[['CYLINDERS','ENGINEIZE','CO2EMISSIONS','FUELCONSUMPTION_COMB']]
vis.hist()
```



```
In [11]: #Plotting these features against the Emission, to see their linear relationship
plt.scatter(cdf.FUELCONSUMPTION_COMB, cdf.CO2EMISSIONS, color='blue')
plt.xlabel("FUELCONSUMPTION_COMB")
plt.ylabel("Emission")
plt.show()
```



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

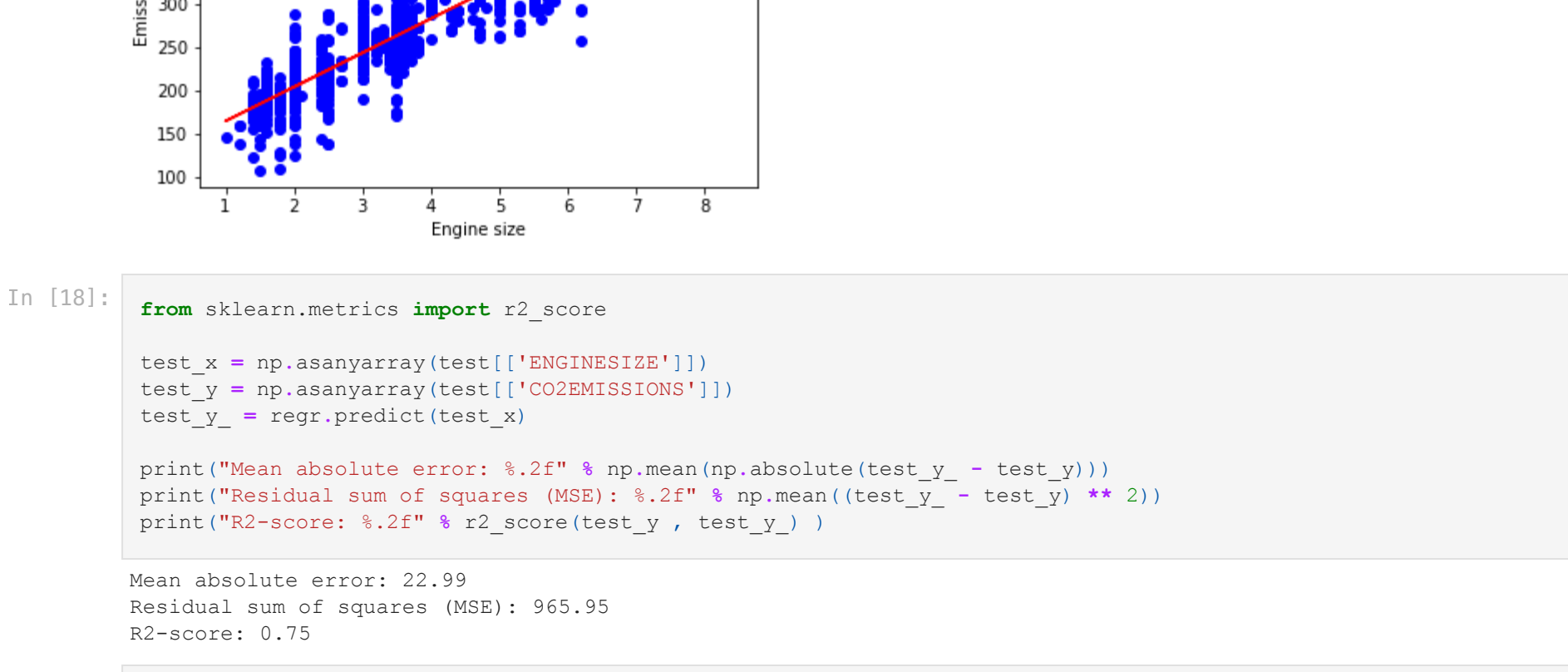


```
In [13]: plt.scatter(cdf.CYLINDERS, cdf.CO2EMISSIONS, color='blue')
plt.xlabel("CYLINDERS")
plt.ylabel("Emission")
plt.show()
```



```
In [14]: #training and testing the data set
msk = np.random.rand(len(df)) < 0.8
train = cdf[msk]
test = cdf[~msk]
```

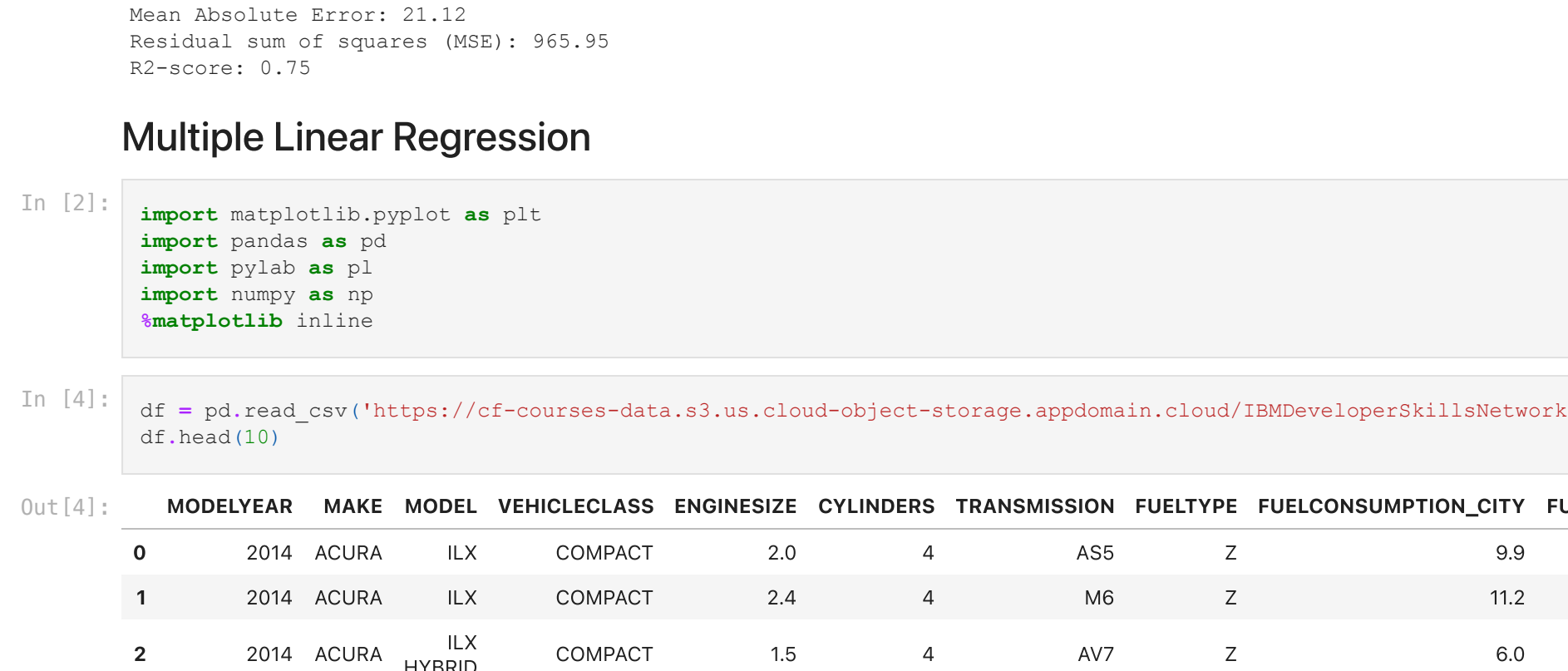
```
In [15]: #For single linear regression model, Linear Regression fits a linear model with coefficients B = (B1, ..., Bn)
#train data distribution
plt.scatter(train.ENGINEIZE, train.CO2EMISSIONS, color='blue')
plt.xlabel("Engine size")
plt.ylabel("Emission")
plt.show()
```



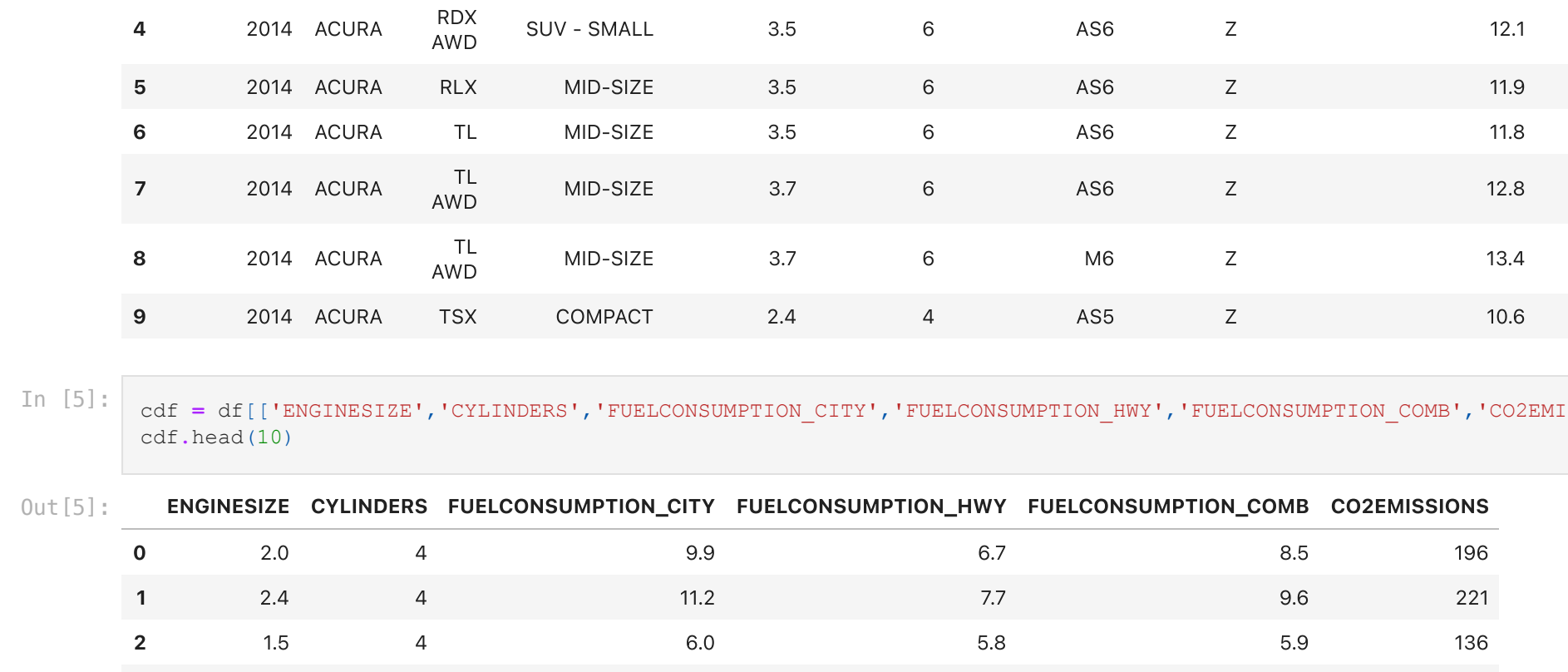
```
In [16]: #Modeling the data using sklearn
from sklearn import linear_model
regr = linear_model.LinearRegression()
train_x = np.asanyarray(train[['ENGINEIZE']])
train_y = np.asanyarray(train[['CO2EMISSIONS']])
regr.fit(train_x, train_y)
# the coefficients
print ('Coefficients: ', regr.coef_)
print ('Intercept: ', regr.intercept_)
```

```
Coefficients: [[39.4855416]]
Intercept: [125.43722767]
```

```
In [17]: #Plotting the fit line over the data
plt.scatter(train.ENGINEIZE, 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")
plt.show()
```



```
Out[17]: Text(0, 0.5, 'Emission')
```



```
In [18]: from sklearn.metrics import r2_score
test_x = np.asanyarray(test[['ENGINEIZE']])
test_y = np.asanyarray(test[['CO2EMISSIONS']])
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_))
```

```
Mean absolute error: 22.99
Residual sum of squares (MSE): 965.95
R2-score: 0.75
```

```
In [19]: #To determine the evaluation metrics using the FUELCONSUMPTION_COMB feature
train_x = train[['FUELCONSUMPTION_COMB']]
regr = linear_model.LinearRegression()
regr.fit(train_x, train_y)
predictions = regr.predict(test_x)
print("Mean absolute error: %.2f" % np.mean(np.absolute(predictions - test_y)))
print("Residual sum of squares (MSE): %.2f" % np.mean((test_y - predictions) ** 2))
print("R2-score: %.2f" % r2_score(test_y, predictions))
```

```
Mean Absolute Error: 21.12
Residual sum of squares (MSE): 965.95
R2-score: 0.75
```

Multiple Linear Regression

```
In [2]: import matplotlib.pyplot as plt
import pandas as pd
import pylab as pl
import numpy as np
%matplotlib inline
```

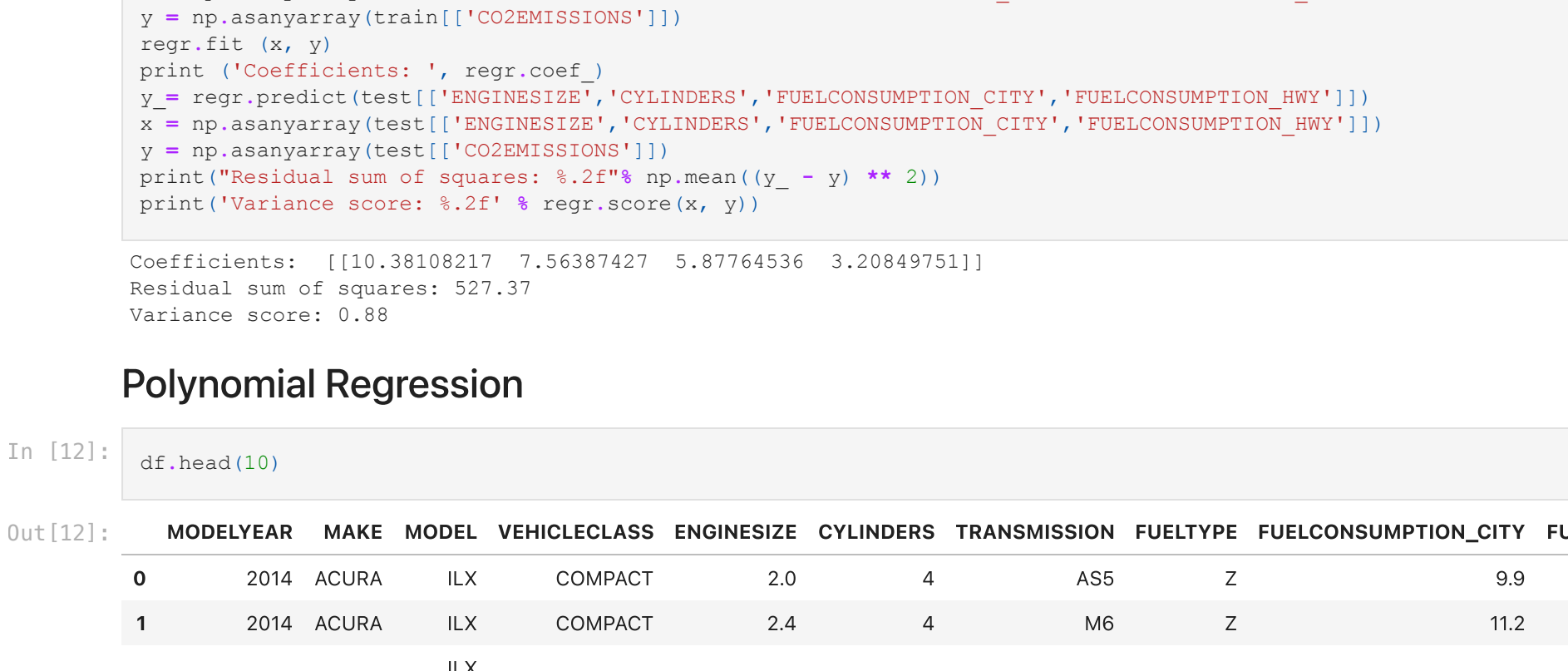
```
In [4]: df = pd.read_csv('https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBMDeveloperSkillsNetwork-
df.head(10)
```

	MODEL	YEAR	MAKE	MODEL	VEHICLECLASS	ENGINEIZE	CYLINDERS	TRANSMISSION	FUELTYPE	FUELCONSUMPTION_CITY	FUELCONSUMPTION_Hwy	CO2EMISSIONS
0	2014	ACURA	ILX	COMPACT	2.0	4	A55	Z		9.9		
1	2014	ACURA	ILX	COMPACT	2.4	4	M6	Z		11.2		
2	2014	ACURA	ILX	HYBRID	COMPACT	1.5	4	AV7	Z	6.0		
3	2014	ACURA	MDX	SUV - SMALL	3.5	6	A56	Z		12.7		
4	2014	ACURA	RDX	SUV - SMALL	3.5	6	A56	Z		12.1		
5	2014	ACURA	RLX	MID-SIZE	3.5	6	A56	Z		11.9		
6	2014	ACURA	TL	MID-SIZE	3.5	6	A56	Z		11.8		
7	2014	ACURA	TL	AWD	MID-SIZE	3.7	6	A56	Z	12.8		
8	2014	ACURA	TL	AWD	MID-SIZE	3.7	6	M6	Z	13.4		
9	2014	ACURA	TSX	COMPACT	2.4	4	A55	Z		10.6		

```
In [5]: cdf = df[['ENGINEIZE','CYLINDERS','FUELCONSUMPTION_CITY','FUELCONSUMPTION_Hwy','FUELCONSUMPTION_COMB','CO2EMISSIONS']]
cdf.head(10)
```

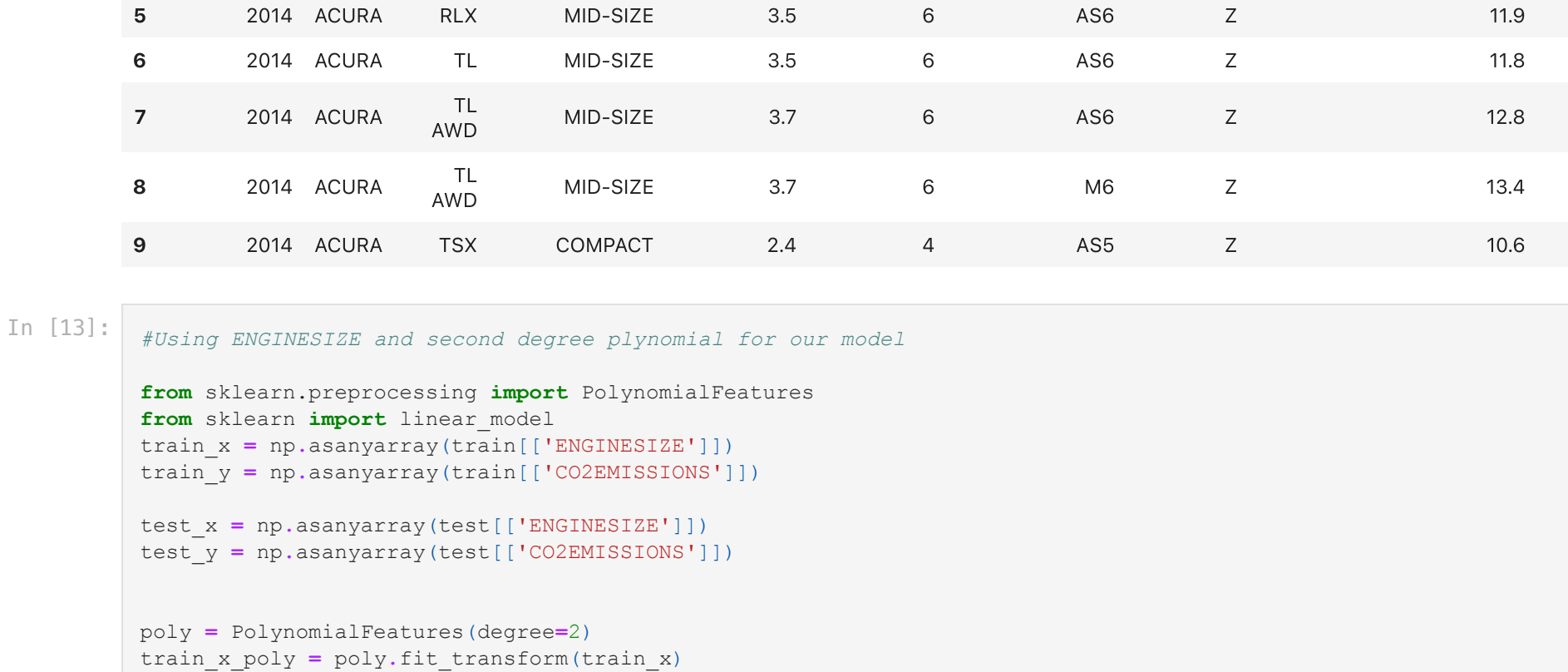
	ENGINEIZE	CYLINDERS	FUELCONSUMPTION_CITY	FUELCONSUMPTION_Hwy	FUELCONSUMPTION_COMB	CO2EMISSIONS
0	2.0	4	9.9	6.7	8.5	196
1	2.4	4	11.2	7.7	9.6	221
2	1.5	4	6.0	5.8	5.9	136
3	3.5	6	12.7	9.1	11.1	255
4	3.5	6	12.1	8.7	10.6	244
5	3.5	6	11.9	7.7	10.0	230
6	3.5	6	11.8	8.1	10.1	232
7	3.7	6	12.8	9.0	11.1	255
8	3.7	6	13.4	9.5	11.6	267
9	2.4	4	10.6	7.5	9.2	212

```
In [6]: plt.scatter(cdf.ENGINEIZE, cdf.CO2EMISSIONS, color='blue')
plt.xlabel("Engine size")
plt.ylabel("Emission")
plt.show()
```



```
In [7]: #Test and training data
msk = np.random.rand(len(df)) < 0.8
train = cdf[msk]
test = cdf[~msk]
```

```
In [8]: plt.scatter(train.ENGINEIZE, train.CO2EMISSIONS, color='blue')
plt.xlabel("Engine size")
plt.ylabel("Emission")
plt.show()
```



```
In [9]: # Using multiple independent variables such as FUELCONSUMPTION_COMB, EngineSize and Cylinders to predict CO2EMISSIONS
#Multiple linear regression model is the extension of the simple linear regression model.
from sklearn import linear_model
regr = linear_model.LinearRegression()
train_x = np.asanyarray(train[['ENGINEIZE','CYLINDERS','FUELCONSUMPTION_COMB']])
train_y = np.asanyarray(train[['CO2EMISSIONS']])
regr.fit(train_x, train_y)
print ('Coefficients: ', regr.coef_)
```

```
Coefficients: [[10.30379996 7.88799292 9.32356035]]
```

```
In [10]: #Estimating the unknown parameters in a linear regression model using Ordinary Least Squares (OLS) method.
y_hat = regr.predict(test[['ENGINEIZE','CYLINDERS','FUELCONSUMPTION_COMB']])
x = np.asanyarray(test[['ENGINEIZE','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))
```

```
Residual sum of squares: 527.94
Variance score: 0.88
```

```
In [11]: # Multiple linear regression using FUELCONSUMPTION_CITY and FUELCONSUMPTION_Hwy instead of FUELCONSUMPTION_COMB
train_x = np.asanyarray(train[['ENGINEIZE','CYLINDERS','FUELCONSUMPTION_CITY','FUELCONSUMPTION_Hwy']])
train_y = np.asanyarray(train[['CO2EMISSIONS']])
regr = linear_model.LinearRegression()
regr.fit(train_x, train_y)
print ('Coefficients: ', regr.coef_)
```

```
Coefficients: [[10.38108217 7.56387427 5.87764536 3.20849751]]
Residual sum of squares: 527.37
Variance score: 0.88
```

Polynomial Regression

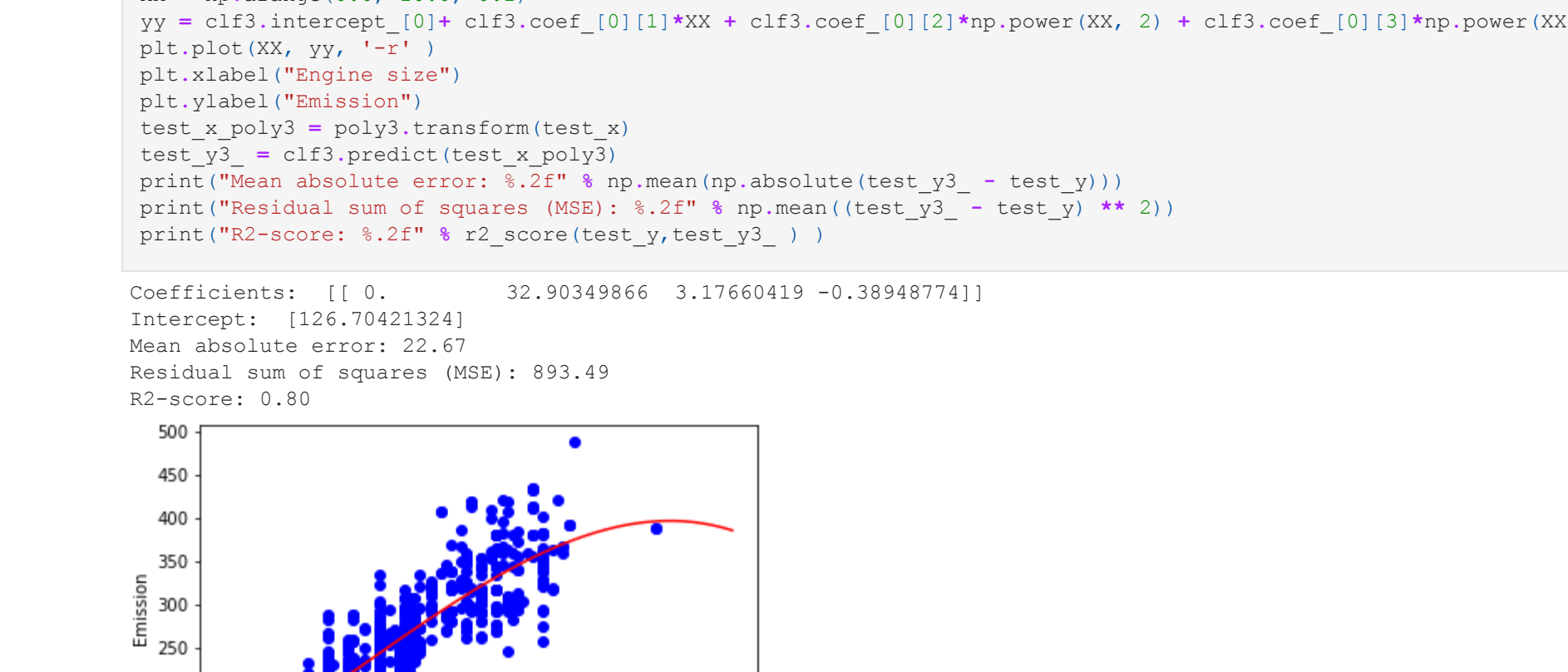
```
In [12]: df.head(10)
```

	MODEL	YEAR	MAKE	MODEL	VEHICLECLASS	ENGINEIZE	CYLINDERS	TRANSMISSION	FUELTYPE	FUELCONSUMPTION_CITY	FUELCONSUMPTION_Hwy	CO2EMISSIONS
0	2014	ACURA	ILX	COMPACT	2.0	4	A55	Z		9.9		
1	2014	ACURA	ILX	COMPACT	2.4	4	M6	Z		11.2		
2	2014	ACURA	ILX	HYBRID	COMPACT	1.5	4	AV7	Z	6.0		
3	2014	ACURA	MDX	SUV - SMALL	3.5	6	A56	Z		12.7		
4	2014	ACURA	RDX	SUV - SMALL	3.5	6	A56	Z		12.1		
5	2014	ACURA	RLX	MID-SIZE	3.5	6	A56	Z		11.9		
6	2014	ACURA	TL	MID-SIZE	3.5	6	A56	Z		11.8		
7	2014	ACURA	TL	AWD	MID-SIZE	3.7	6	A56	Z	12.8		
8	2014	ACURA	TL	AWD	MID-SIZE	3.7	6	M6	Z	13.4		
9	2014	ACURA	TSX	COMPACT	2.4	4	A55	Z		10.6		

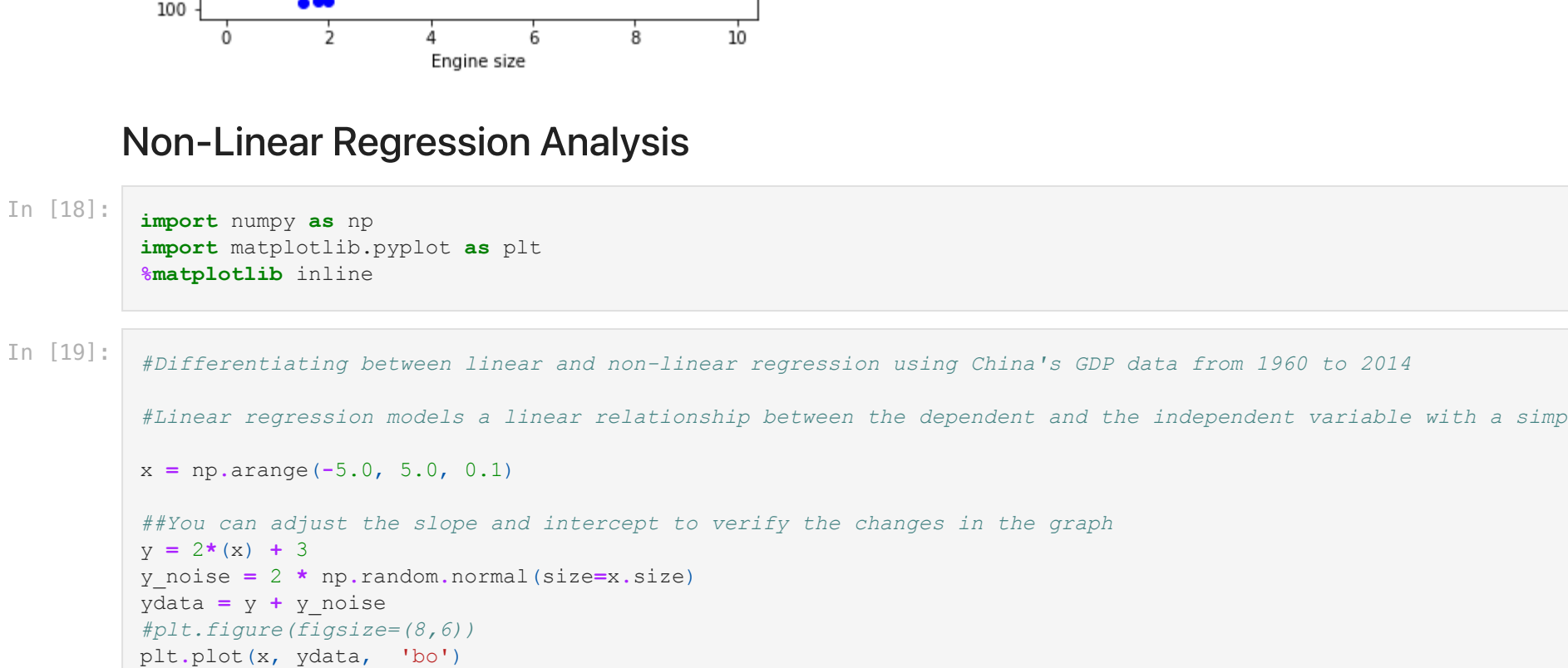
```
In [13]: #Using ENGINEIZE and second degree polynomial for our model
from sklearn.preprocessing import PolynomialFeatures
from sklearn import linear_model
train_x = np.asanyarray(train[['ENGINEIZE','CYLINDERS','FUELCONSUMPTION_COMB']])
train_y = np.asanyarray(train[['CO2EMISSIONS']])
regr = linear_model.LinearRegression()
regr.fit(train_x, train_y)
print ('Coefficients: ', regr.coef_)
```

```
Coefficients: [[0. 0. 50.08730154 -1.54004369]]
Intercept: [108.38613302]
```

```
In [15]: plt.scatter(train.ENGINEIZE, train.CO2EMISSIONS, color='blue')
xx = np.arange(0.0, 10.0, 0.1)
yy = regr.predict([0]*4+xx)
plt.plot(xx, yy, '-r')
plt.xlabel("Engine size")
plt.ylabel("Emission")
plt.show()
```



```
Out[15]: Text(0, 0.5, 'Emission')
```



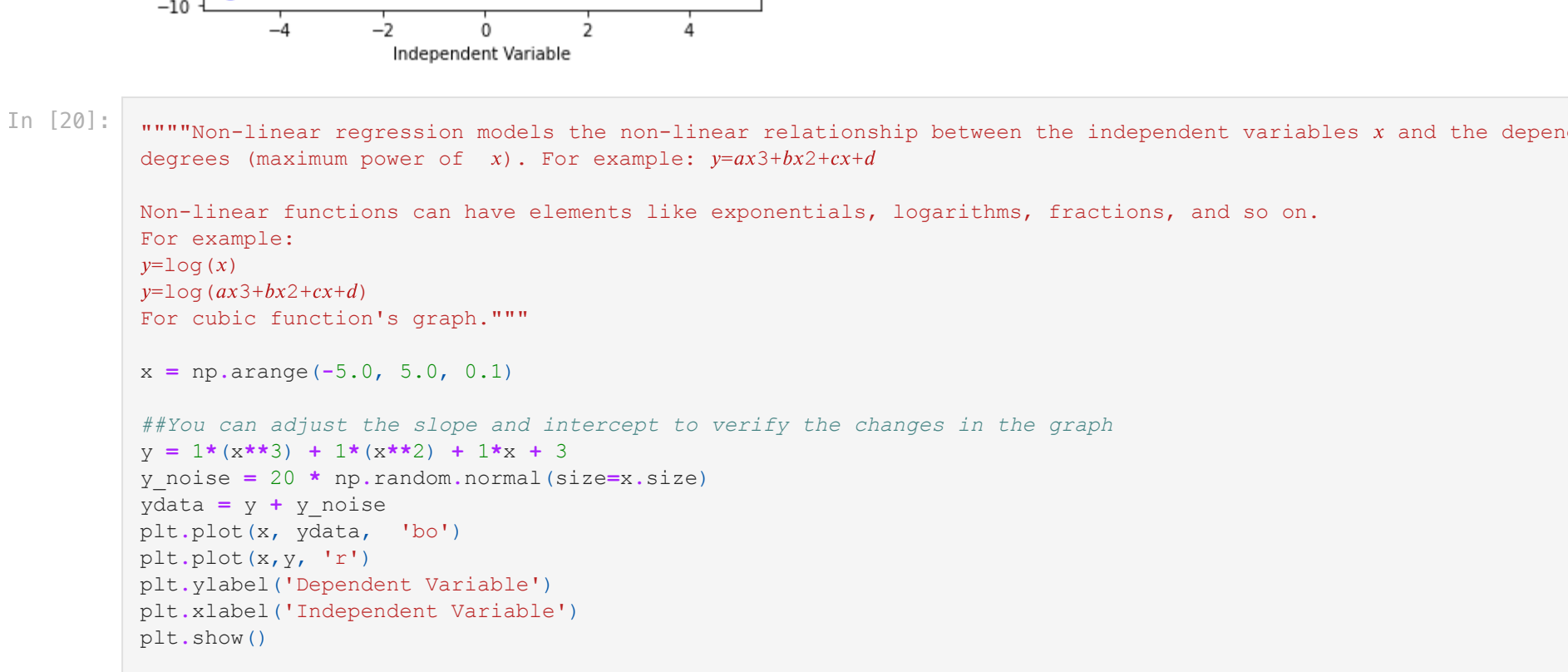
```
In [16]: from sklearn.metrics import r2_score
test_x_poly = poly3.fit_transform(test_x)
test_y_ = regr.predict(test_x_poly)
print("Mean absolute error: %.2f" % np.mean(np.absolute(test_y_ - test_y)))
print("Residual sum of squares (MSE): %.2f" % np.mean((test_y_ - test_y) ** 2))
print("R2-score: %.2f" % r2_score(test_y_, test_y_))
```

```
Mean absolute error: 22.76
Residual sum of squares (MSE): 901.24
R2-score: 0.80
```

```
In [17]: # Polynomial regression using above data with degree of three (cubic)
poly3 = PolynomialFeatures(degree=3)
train_x_poly3 = poly3.fit_transform(train_x)
clf3 = linear_model.LinearRegression()
train_y_ = clf3.fit(train_x_poly3, train_y)
```

```
print ('Coefficients: ', clf3.coef_)
print ('Intercept: ', clf3.intercept_)
plt.scatter(train.ENGINEIZE, train.CO2EMISSIONS, color='blue')
xx = np.arange(0.0, 10.0, 0.1)
yy = clf3.predict([0]*4+xx)
plt.plot(xx, yy, '-r')
plt.xlabel("Engine size")
plt.ylabel("Emission")
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_))
```

```
Coefficients: [[0. 0. 32.90349866 3.17660419 -0.38948774]]
Intercept: [126.70421324]
Mean absolute error: 22.67
Residual sum of squares (MSE): 893.49
R2-score: 0.80
```



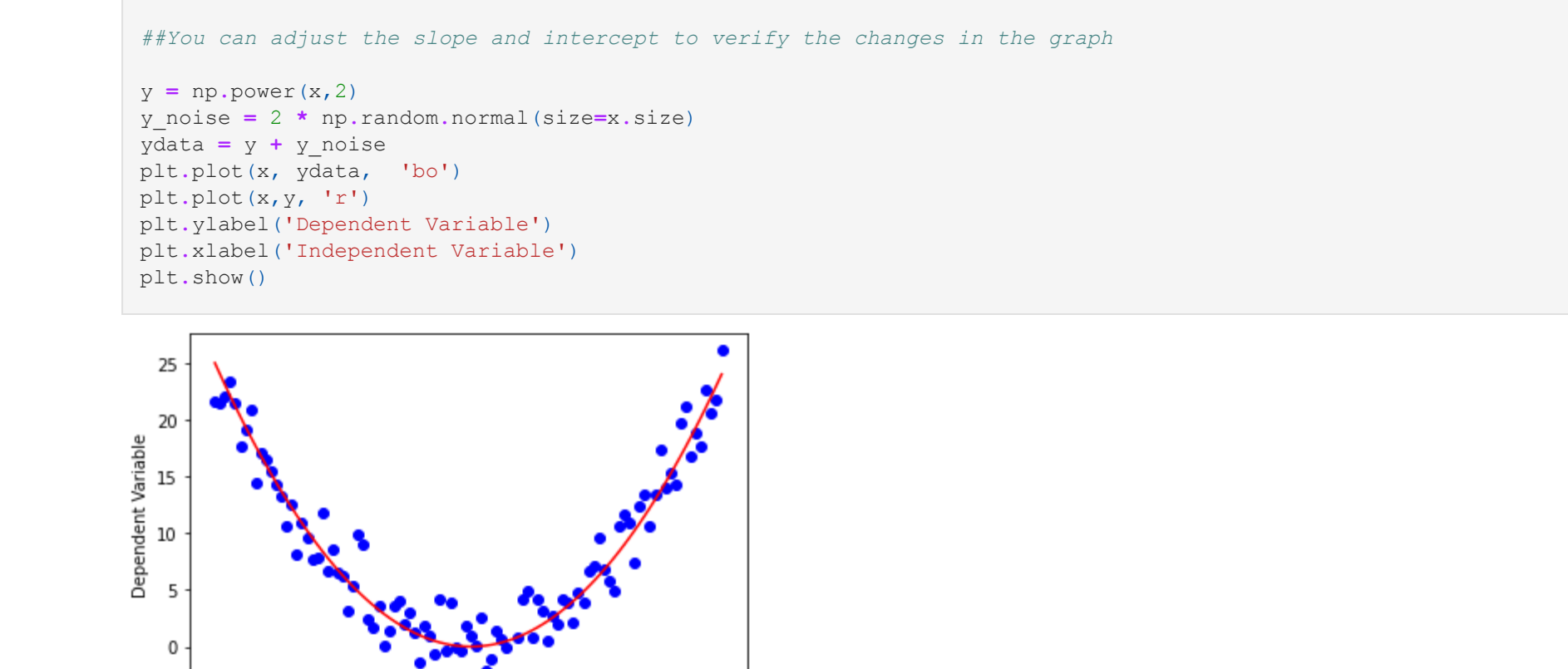
Non-Linear Regression Analysis

```
In [18]: import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
```

```
In [19]: #Differentiating between linear and non-linear regression using China's GDP data from 1960 to 2014
#linear regression models a linear relationship between the dependent and the independent variable with a slope.
x = np.arange(-5.0, 5.0, 0.1)
```

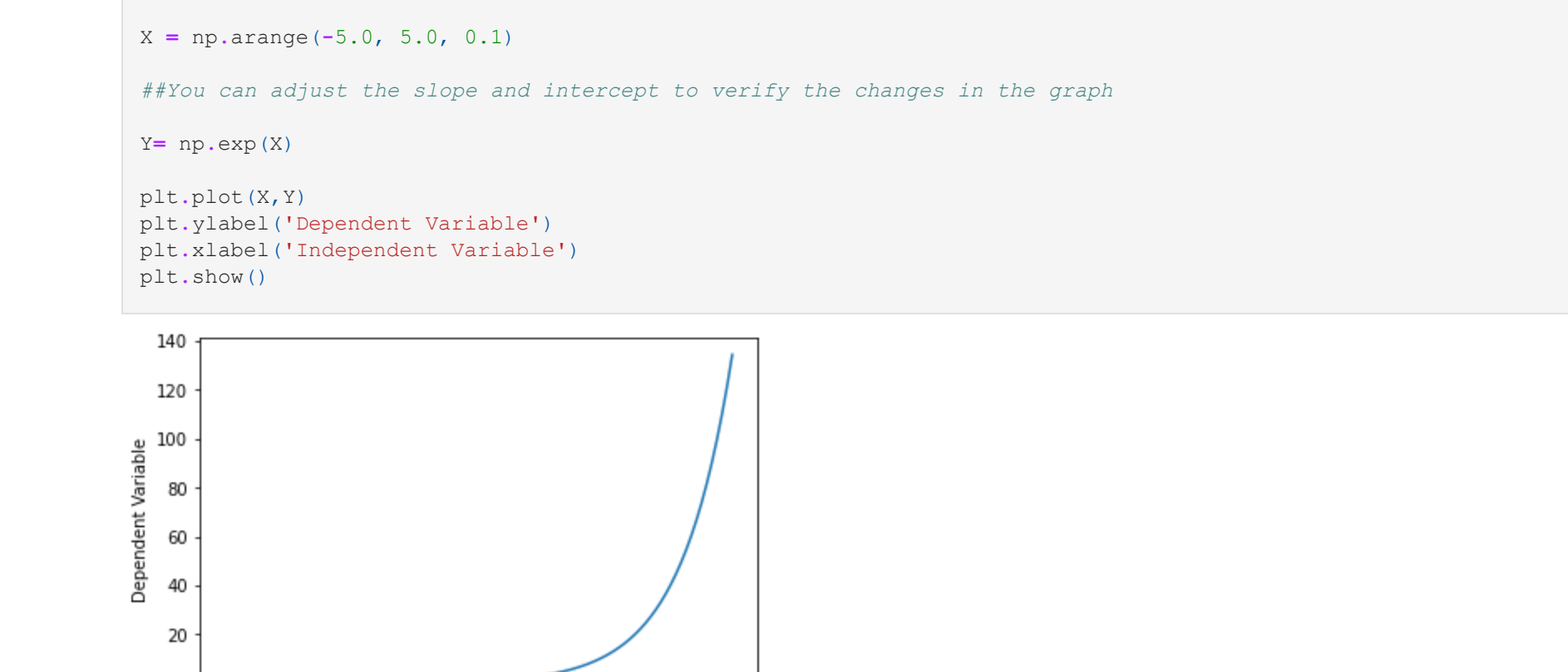
```
##You can adjust the slope and intercept to verify the changes in the graph
y = 2*(x+3) + 10*(x**2) + 1*(x**3)
y_noise = 2 * np.random.normal(size=x.size)
ydata = y + y_noise
plt.figure(figsize=(8,6))
plt.plot(x, ydata, 'bo')
plt.plot(x, y, 'r')
```

```
plt.ylabel('Dependent Variable')
plt.xlabel('Independent Variable')
plt.show()
```



```
In [20]: #Exponential: An exponential function with base c is defined by Y=c*b^X
where b > 0, c > 0, c != 1, and x is any real number. The base, c, is constant and the exponent, x, is a variable.
In general form it would be written as y=log(X)
X = np.arange(-5.0, 5.0, 0.1)
Y = np.log(X)
```

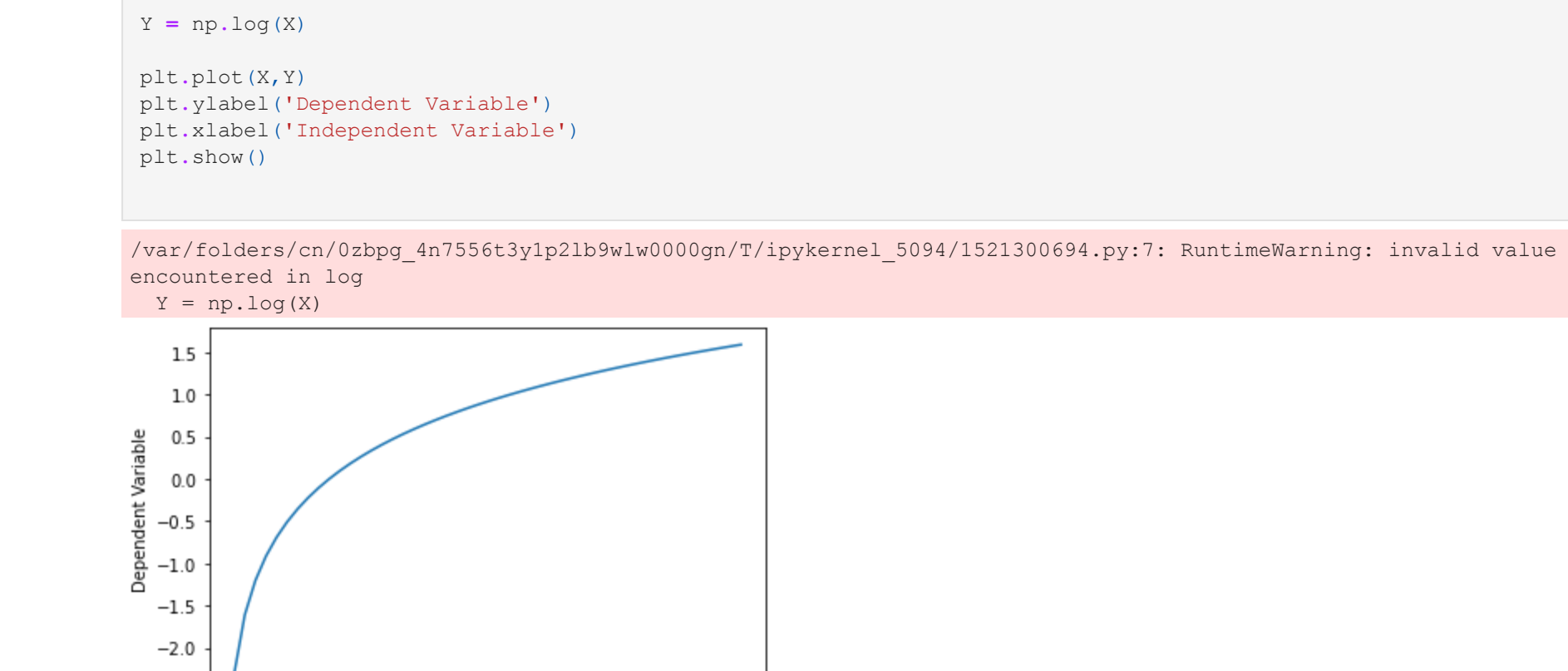
```
plt.plot(X,Y)
plt.ylabel('Dependent Variable')
plt.xlabel('Independent Variable')
plt.show()
```



```
In [21]: #For quadratic function's graph (Y=X^2(squared))
x = np.arange(-5.0, 5.0, 0.1)
```

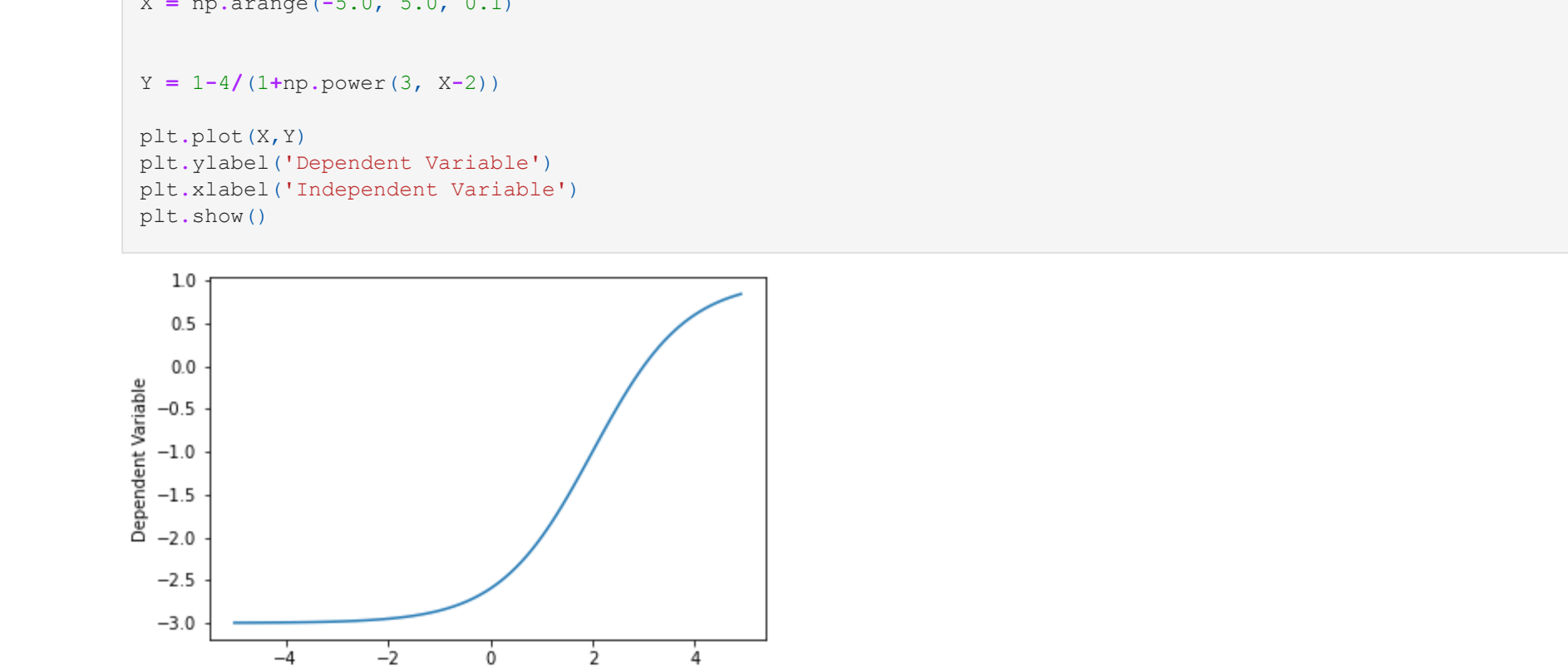
```
##You can adjust the slope and intercept to verify the changes in the graph
y = np.power(x,2)
y_noise = 2 * np.random.normal(size=x.size)
ydata = y + y_noise
plt.plot(x, ydata, 'bo')
plt.plot(x, y, 'r')
```

```
plt.ylabel('Dependent Variable')
plt.xlabel('Independent Variable')
plt.show()
```



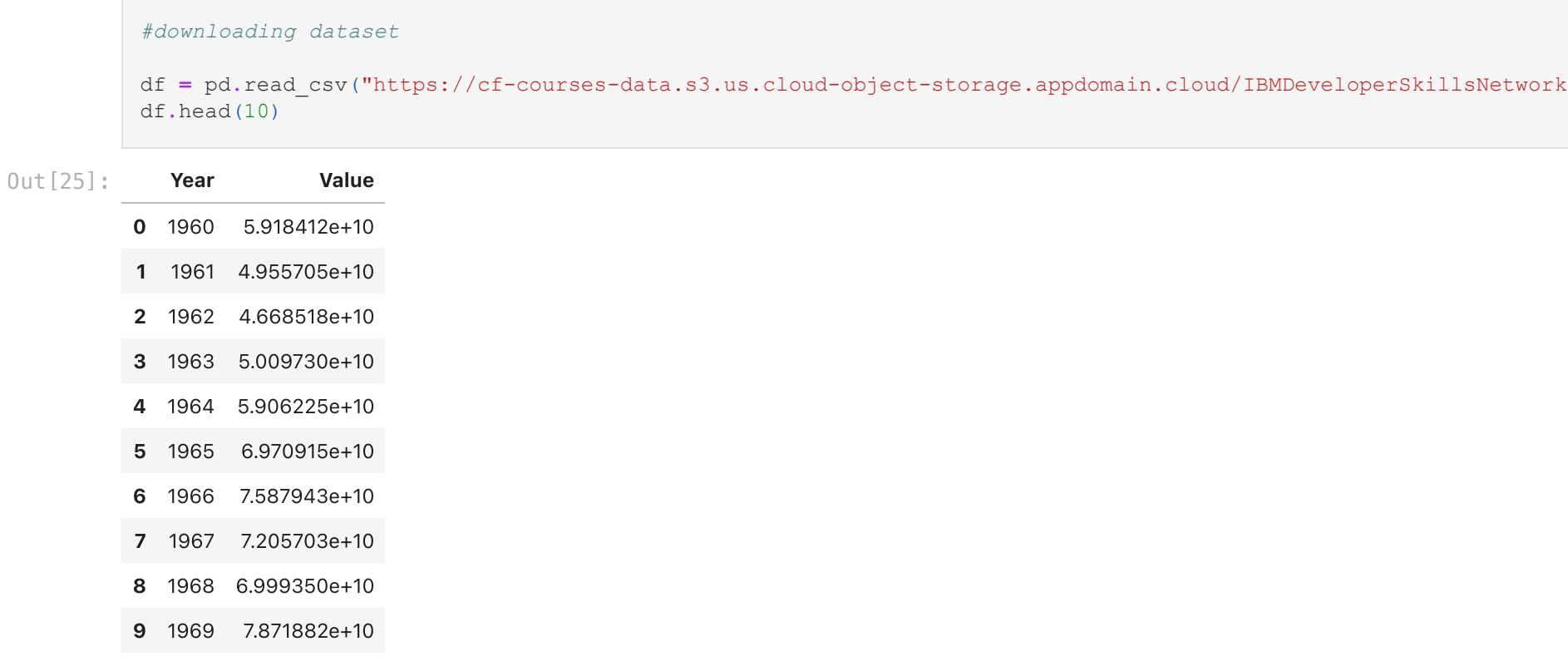
```
In [22]: #Exponential: An exponential function with base c is defined by Y=c*b^X
where b > 0, c > 0, c != 1, and x is any real number. The base, c, is constant and the exponent, x, is a variable.
In general form it would be written as y=log(X)
X = np.arange(-5.0, 5.0, 0.1)
Y = np.log(X)
```

```
plt.plot(X,Y)
plt.ylabel('Dependent Variable')
plt.xlabel('Independent Variable')
plt.show()
```



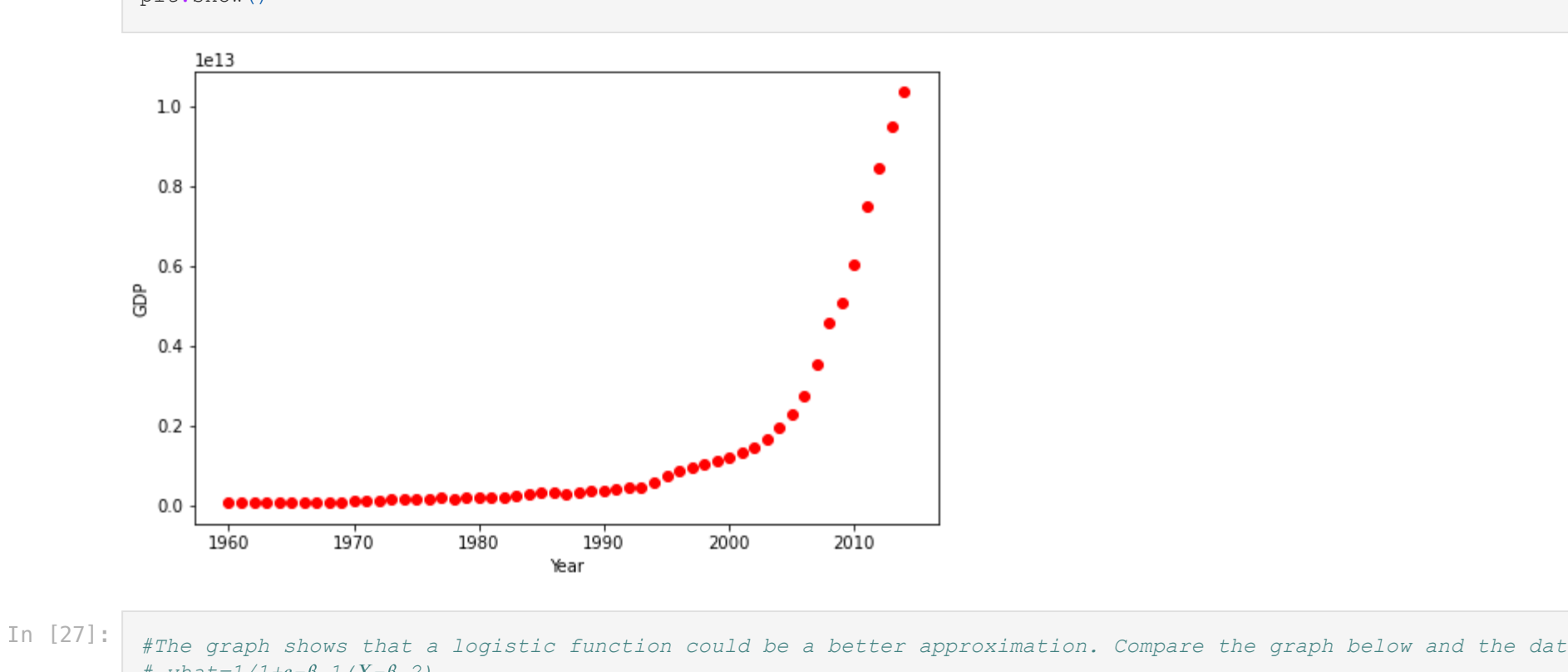
```
In [23]: #Logistic: The response y is a result of applying the logarithmic map from the input x to the output y.
where b > 0, c > 0, c != 1, and x is any real number. The base, c, is constant and the exponent, x, is a variable.
In general form it would be written as y=log(X)
X = np.arange(-5.0, 5.0, 0.1)
Y = np.log(X)
```

```
plt.plot(X,Y)
plt.ylabel('Dependent Variable')
plt.xlabel('Independent Variable')
plt.show()
```



```
In [24]: #Sigmoidal/Logistic: Y = a + b/(1+(X^d)) - (note: X^d is the power of d)
X = np.arange(-5.0, 5.0, 0.1)
Y = 1/(1+np.power(3, X-2))
```

```
plt.plot(X,Y)
plt.ylabel('Dependent Variable')
plt.xlabel('Independent Variable')
plt.show()
```



```
In [25]: #Non-Linear Regression Example
import numpy as np
import pandas as pd
#downloading dataset
df = pd.read_csv('https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBMDeveloperSkillsNetwork-
df.head(10)
```

	Year	Value
0	1960	5.918412e+10
1	1961	4.955705e+10
2	1962	4.668518e+10
3	1963	5.009730e+10
4	1964	5.906225e+10
5	1965	6.970915e+10
6	1966	7.587943e+10
7	1967	7.205703e+10
8	1968	6.999350e+10
9	1969	7.871882e+10

```
In [26]: plt.figure(figsize=(8,5))
x_data, y_data = (df["Year"].values, df["Value"].values)
plt.plot(x_data, y_data, 'ro')
plt.xlabel("Year")
plt.ylabel("Value")
plt.show()
```



```
In [27]: #The graph shows that a logistic function could be a better approximation. Compare the graph below and the data.
figsize=(14,8)
X = np.arange(-5.0, 5.0, 0.1)
Y = 1/(1.0 + np.exp(-X))
```

```
plt.plot(X,Y)
plt.ylabel('Dependent Variable')
plt.xlabel('Independent Variable')
plt.show()
```



```
[28]: #Building the mode
def sigmoid(x, Beta_1, Beta_2):
    y = 1 / (1 + np.exp(-Beta_1*(x-Beta_2)))
    return y

In [29]: beta_1 = 0.10
beta_2 = 1990.0

#sigmoid function
Y_pred = sigmoid(x_data, beta_1 , beta_2)

#plot initial prediction against datapoints
plt.plot(x_data, Y_pred*1500000000000.0)
plt.plot(x_data, y_data, 'ro')

Out[29]:
<matplotlib.lines.Line2D at 0x7ffe6c781340>

In [30]: # normalizing the data
xdata = x_data/max(x_data)
ydata = y_data/max(y_data)

In [31]: from scipy.optimize import curve_fit
popt, pcov = curve_fit(sigmoid, xdata, ydata)
#print the final parameters
print(" beta_1 = %f, beta_2 = %f" % (popt[0], popt[1]))

beta_1 = 690.453017, beta_2 = 0.997207

In [32]: #Resulting regression
x = np.linspace(1960, 2015, 55)
x = x/max(x)
plt.figure(figsize=(8,5))
y = sigmoid(x, *popt)
plt.plot(xdata, ydata, 'ro', label='data')
plt.plot(x,y, linewidth=3.0, label='fit')
plt.legend(loc='best')
plt.ylabel('GDP')
plt.xlabel('Year')
plt.show()

In [33]: #Calculating the accuracy of the model
# split data into train/test
msk = np.random.rand(len(df)) < 0.8
train_x = xdata[msk]
test_x = xdata[~msk]
train_y = ydata[msk]
test_y = ydata[~msk]

# build the model using train set
popt, pcov = curve_fit(sigmoid, train_x, train_y)

# predict using test set
y_hat = sigmoid(test_x, *popt)

# evaluation
print("Mean absolute error: %.2f" % np.mean(np.absolute(y_hat - test_y)))
print("Residual sum of squares (MSE): %.2f" % np.mean((y_hat - test_y) ** 2))
from sklearn.metrics import r2_score
print("R2-score: %.2f" % r2_score(test_y,y_hat) )

Mean absolute error: 0.04
Residual sum of squares (MSE): 0.00
R2-score: 0.96
```