



VISVESVARAYA NATIONAL INSTITUTE
OF TECHNOLOGY (VNIT), NAGPUR

Machine Learning with Python (ECL 443)

Lab Report

Submitted by :

Pranay Dumbhare (BT19ECE089)

Semester 7

Submitted to :

Dr.Saugata Sinha

(Course Instructor)

Department of Electronics and Communication Engineering

VNIT Nagpur

Contents

1	Linear Regression	2
1.1	Abstract	2
1.2	Introduction	2
1.3	Problem Statement	2
1.4	Method/Procedure	2
1.5	Results/Discussion	4
1.6	Conclusion	6
1.7	Appendices	6

Linear Regression

1.1 Abstract: Machine Learning is a computer program which is used to train our machine to predict some outputs on the basis of previous inputs. With advancements in machine learning there are many algorithms available to find patterns in the given data so that the machine can learn from these patterns to predict future output. The accuracy of the model is how correctly it predicts the output.

1.2 Introduction: Linear Regression is a machine learning algorithm where there is a linear dependency between the dependent variable (y) and independent variables (x1,x2,...,xn). Following is the general representation of linear regression:

$$y = w_0 + w_1 * x_1 + w_2 * x_2 + \dots + w_n * x_n$$

Where, y is the dependent variable,
x1,x2,...,xn are the independent variables and
w0,w1,...wn are the weights of these dependent variables which resembles how much the dependent variable(y) depends on the independent variable(xn).

1.3 Problem Statement: Predict the number of traffic fatalities caused on the basis of number of registered vehicles (in thousands), number of licensed drivers (in thousands) and the number of miles travelled (in millions) by the vehicles.

1.4 Method/Procedure:

Note: The given matlab file (accidents.mat) contains many columns(both dependent and independent). From these columns, I have chosen the number of traffic fatalities as dependent variable and number of registered vehicles, number of licensed drivers and the number of miles travelled by the vehicles as independent variables.

1. Load the given data file is **Matlab_accidents.mat** using the user-defined function **BT19ECE089_train_test_split**. The data is loaded as a pandas dataframe.
2. After the data is loaded into the dataframe, a new dataframe is created using the selected dependent and independent variables.
3. Normalize the data frame column wise.
4. The data is then split into train and test with split ratio 0.3. Then matrix theta train , y train , theta test , y test are created to perform training and testing.

5. Pseudo Inverse Method

- (a) Compute optimum weights by computing pseudo inverse of the theta matrix and multiplying it with y train.
- (b) Predict the output (y pred) by multiplying this optimum weights with theta train.
- (c) Plot y pred and y test to visualize how accurately our model predicts the output.
- (d) Calculate root mean square error, mean squared error and mean absolute error to analyze model analytically.

6. Gradient Descent Method

- (a) Compute optimum weights using gradient descent algorithm. Here we use number of iterations as 1000 and learning rate as 0.03
 - (b) Predict the output (y pred) by multiplying this optimum weights with theta train.
 - (c) Plot y pred and y test to visualize how accurately our model predicts the output.
 - (d) Calculate root mean square error, mean squared error and mean absolute error to analyze model analytically.
7. For change in relationship between the input and output variables, we consider that the fatalities (y) varies as square of Registered vehicles (thousands) (x1). This term is added after squaring x1, and the accuracy is calculated again.

1.5 Results/Discussion:

Note: The given matlab file (*accidents.mat*) contains many columns (both dependent and independent). From this columns, I have chosen the number of traffic fatalities as dependent variable and number of registered vehicles, number of licensed drivers and the number of miles travelled by the vehicles as independent variables.

1. **Pseudo Inverse Method** The RMS Accuracy obtained using the Pseudo Inverse Method is 96%.

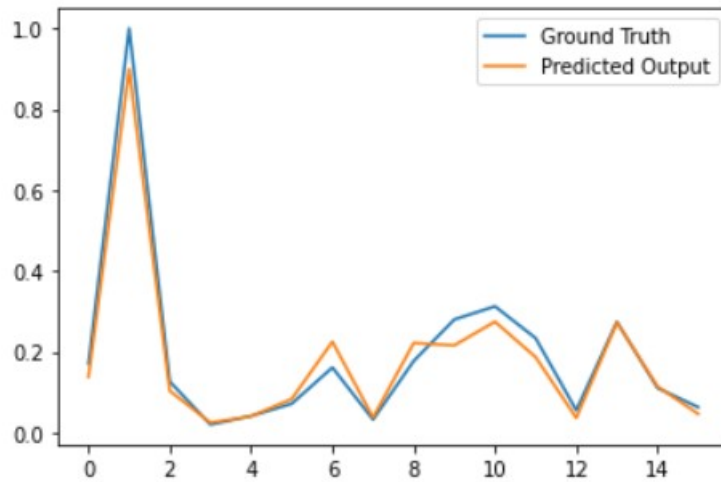


Figure 1: Predicted Output and Ground Truth using Pseudo Inverse

Mean Squared Error: 0.0016484915795305277
 Root Mean Squared Error: 0.04060162040523171
 Mean Absolute Error: 0.0297834829376438

2. **Gradient Descent Method** The RMS obtained is 94.6%. As we can see, this is a significant improvement over the Pseudo Inverse method.

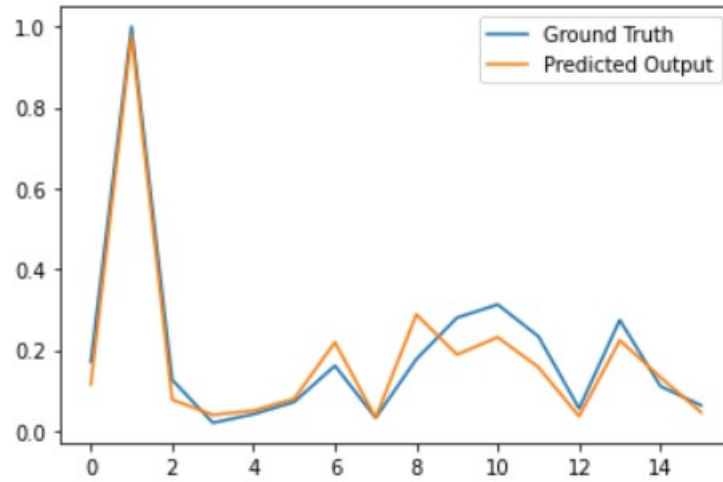


Figure 2: Predicted Output and Ground Truth using Gradient Descent

Mean Squared Error: 0.002906374104840904
 Root Mean Squared Error: 0.05391079766466922
 Mean Absolute Error: 0.043063532175202424

- For change in relationship between the input and output variables, we consider that the fatalities (y) varies as square of Registered vehicles (thousands) (x1). This term is added after squaring x1. The RMS Accuracy obtained is 95.4%.

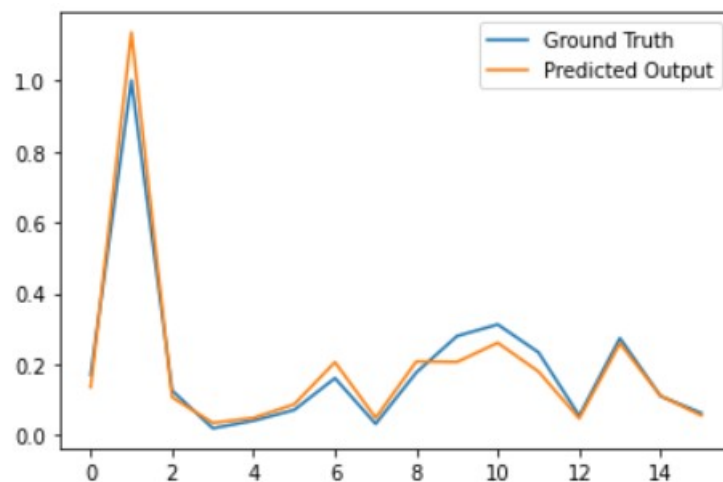


Figure 3: Predicted Output and Ground Truth After changing the relationship between input and output

```

Mean Squared Error: 0.0021891078643436262
Root Mean Squared Error: 0.04678790297014418
Mean Absolute Error: 0.033153048552636656

```

1.6 Conclusion: Linear Regression formulates a linear relationship between the independent variables as functions of dependent variables with some weights associated with each dependent variable which resembles how much the dependent variable(y) depends on the independent variable(xn). The model then tries to find the optimum value of these weights so as to minimize the error or loss.

This finding of optimum weights can be done by using two methods:

- Pseudo Inverse
- Gradient Descent

The accuracy of pseudo inverse is better than the gradient descent as gradient descent is approximation based algorithm while pseudo inverse calculates the weights precisely.

1.7 Appendices: The code for linear regression is given below:

```

1  #Installing Necessary Libraries
2
3  !pip install mat4py
4  import pandas as pd
5  import numpy as np
6  from sklearn.model_selection import train_test_split
7  from mat4py import loadmat
8  import matplotlib.pyplot as plt
9  from sklearn.metrics import mean_squared_error, mean_absolute_error
10
11 #Function to Read Data File
12
13 def BT19ECE089_dataset_div_shuffle(dir_path, split_ratio):
14     #CSV XLS MAT
15     if(dir_path.endswith('.csv')):
16         data = pd.read_csv(dir_path)
17     elif(dir_path.endswith('.xls')):
18         data = pd.read_excel(dir_path)
19     elif(dir_path.endswith('.mat')):
20         data = loadmat(dir_path)
21     else:
22         print("Invalid datatype. Select .csv or .xls or .mat file")
23

```

```

24     train_data , test_data = train_test_split(data,train_size = ...
        split_ratio,shuffle = True)
25
26     return train_data , test_data
27
28 #Arranging the data frame to get dependent(y) and independent ...
    variables(x)
29
30 data = loadmat('/content/Matlab_accidents.mat')
31
32 data_accidents = data['accidents']
33 #datax = data_accidents[0][2]
34 'statelabel', 'hwydata', 'hwyheaders'
35 df = pd.DataFrame(data_accidents['hwydata'],columns = ...
    data_accidents['hwyheaders'])
36 states = [y[0] for y in data_accidents["statelabel"]]
37 df.insert(loc = 0,column = 'States',value = states)
38 # df.head()
39 df = pd.DataFrame([df['Traffic fatalities'],df['Licensed drivers ...
    (thousands)'],df['Registered vehicles ...
    (thousands)'],df['Vehicle-miles traveled ...
    (millions)']]).transpose()
40
41 #Normalizing and Splitting the data into train and test
42
43 for column in df:
44     df[column] = df[column]/np.amax(df[column])
45 train , test = train_test_split(df,train_size = 0.7,shuffle = True)
46
47 #Creating theta matrix and y for training and testing
48
49 y_train = np.array(train.iloc[:, 0])
50 y_train = y_train.reshape([35,1])
51 x_train = np.array(train.iloc[:, 1:])
52 on_train = np.ones([train.shape[0],1])
53 theta_train = np.hstack((on_train,x_train))
54
55 y_test = np.array(test.iloc[:, 0])
56 x_test = np.array(test.iloc[:, 1:])
57 # x_test = x_test / np.amax(x_test)
58 # y_test = y_test / np.amax(y_test)
59 on_test = np.ones([test.shape[0],1])
60 y_test = y_test.reshape([y_test.shape[0],1])
61 theta_test = np.hstack((on_test,x_test))
62
63 #Calculating optimum weights (using pseudo inverse) and ...
    predicting output
64
65 weight_opt = np.matmul(np.linalg.pinv(theta_train),y_train)

```



```
66 # weight_opt.shape
67 y_pred = np.matmul(theta_test, weight_opt)
68
69 #Plotting ground truth and predicted outputs
70
71 plt.plot(y_test)
72 plt.plot(y_pred)
73
74 #Calculating mean squared error , root mean squared error and ...
    mean absolute error
75
76 mse = mean_squared_error(y_test, y_pred)
77 rmse = mean_squared_error(y_test, y_pred, squared = False)
78 mae = mean_absolute_error(y_test, y_pred)
79 print(mse, rmse, mae)
80
81 #Finding optimum weights using Gradient Descent
82
83 #Randomly initialising weights
84
85 grad_weights = np.random.randn(len(theta_train[0]), 1)
86 # grad_weights.shape
87
88 #Implementing Gradient Descent and predicting output
89
90 iterations = 1000
91 learning_rate = 0.03
92
93 for i in range(iterations):
94     Δ_e = ...
        np.matmul(theta_train.transpose(), (np.matmul(theta_train, grad_weights) - y_train))
95     grad_weights = grad_weights - learning_rate * Δ_e
96 ygrad_pred = np.matmul(theta_test, grad_weights)
97
98 #Plotting ground truth and predicted outputs
99
100 plt.plot(y_test)
101 plt.plot(ygrad_pred)
102
103 #Calculating mean squared error , root mean squared error and ...
    mean absolute error
104
105 mse_g = mean_squared_error(y_test, ygrad_pred)
106 rmse_g = mean_squared_error(y_test, ygrad_pred, squared = False)
107 mae_g = mean_absolute_error(y_test, ygrad_pred)
108 print(mse_g, rmse_g, mae_g)
109
110 #Changing relationship between input and output variables
111
```

```
112 #Recalculating Theta matrix
113
114 theta_train_dash = ...
      np.hstack((on_train,np.reshape(np.square(x_train.T[0]).T,(35,1)),x_train))
115 theta_test_dash = ...
      np.hstack((on_test,np.reshape(np.square(x_test.T[0]).T,(16,1)),x_test))
116 # theta_train_dash.shape
117
118 #Calculating optimum weights (using pseudo inverse) and ...
      predicting output
119
120 weight_opt_dash = ...
      np.matmul(np.linalg.pinv(theta_train_dash),y_train)
121 # weight_opt.shape
122 y_pred_dash = np.matmul(theta_test_dash,weight_opt_dash)
123
124 #Plotting ground truth and predicted outputs
125
126 plt.plot(y_test)
127 plt.plot(y_pred_dash)
128
129 #Calculating mean squared error , root mean squared error and ...
      mean absolute error
130
131 mse_dash = mean_squared_error(y_test, y_pred_dash)
132 rmse_dash = mean_squared_error(y_test, y_pred_dash, squared = False)
133 mae_dash = mean_absolute_error(y_test, y_pred_dash)
134 print(mse_dash,rmse_dash,mae_dash)
```