

IE406 Machine Learning

Lab Assignment - 1

Group 14

201901466: Miti Purohit

202001430: Aryan Shah

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from mpl_toolkits import mplot3d
4 import pandas as pd
```

Listing 1: Libraries

Question 1

Plot θ vs. $L(\theta)$, where $L(\theta) = \theta^2$. θ varies from -10 to +10 with a step size of 0.1. Now locate the minimum value of $L(\theta)$ with corresponding θ value from the plot.

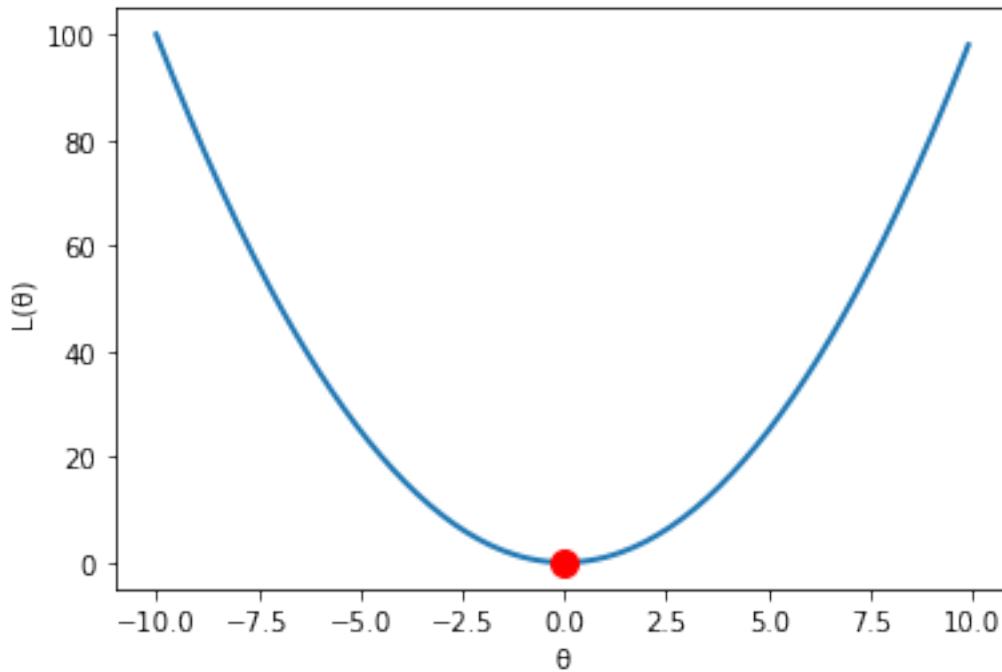
Answer

code

```
1 def L(theta):
2     return theta**2
3
4 theta = np.arange(-10, 10, 0.1)
5 func_L = L(theta)
6 plt.figure()
7 plt.plot(theta, func_L, linewidth=2)
8
9 min_func_L = min(func_L)
10 min_location = np.argwhere(func_L == min_func_L)
11 #print(min_location)
12 plt.plot(theta[min_location[0,0]], min_func_L, 'ro', markersize=10)
13 plt.xlabel(' ')
14 plt.ylabel('L( )')
15 plt.show()
```

Listing 2: Question 1

Result



Observation/ Justification

The minimum value of $L(\theta)$ is 0.0 at 0.0

Question 2

Plot θ vs. $L(\theta)$, where $L(\theta) = \theta_1^2 + \theta_2^2$. θ_1 and θ_2 vary from -10 to +10 with a step size of 0.1. Locate the minimum value of $L(\theta_1, \theta_2)$ with corresponding θ_1, θ_2 values from the plot.

Answer

code

```
1 ddef L(theta1, theta2):
2     return theta1**2 + theta2**2
3
4 theta1 = np.arange(-10,10,0.1)
5 theta2 = np.arange(-10,10,0.1)
6
7 func_L = np.zeros((len(theta1),len(theta2)))
8 for i in range(len(theta1)):
9     for j in range(len(theta2)):
10        func_L[i][j] = L(theta1[i], theta2[j])
11
12 X, Y = np.meshgrid(theta1, theta2)
13 z = L(X, Y)
14 fig = plt.figure(figsize=(10, 10))
15 ax = plt.axes(projection='3d')
16 ax.plot_wireframe(X, Y, z, linewidth=0.3)
17
18
19 min_func_L = np.min(func_L)
20 #print(min_func_L)
21 min_loc = np.argwhere(func_L == min_func_L)
22 #print(z[100,100])
23 x_coord = min_loc[0,0]
24 y_coord = min_loc[0,1]
25 min_value = z[x_coord, y_coord]
26 #print(min_value)
27 ax.scatter(theta1[x_coord], theta2[y_coord], min_value, marker="o", c="red", s=28)
```

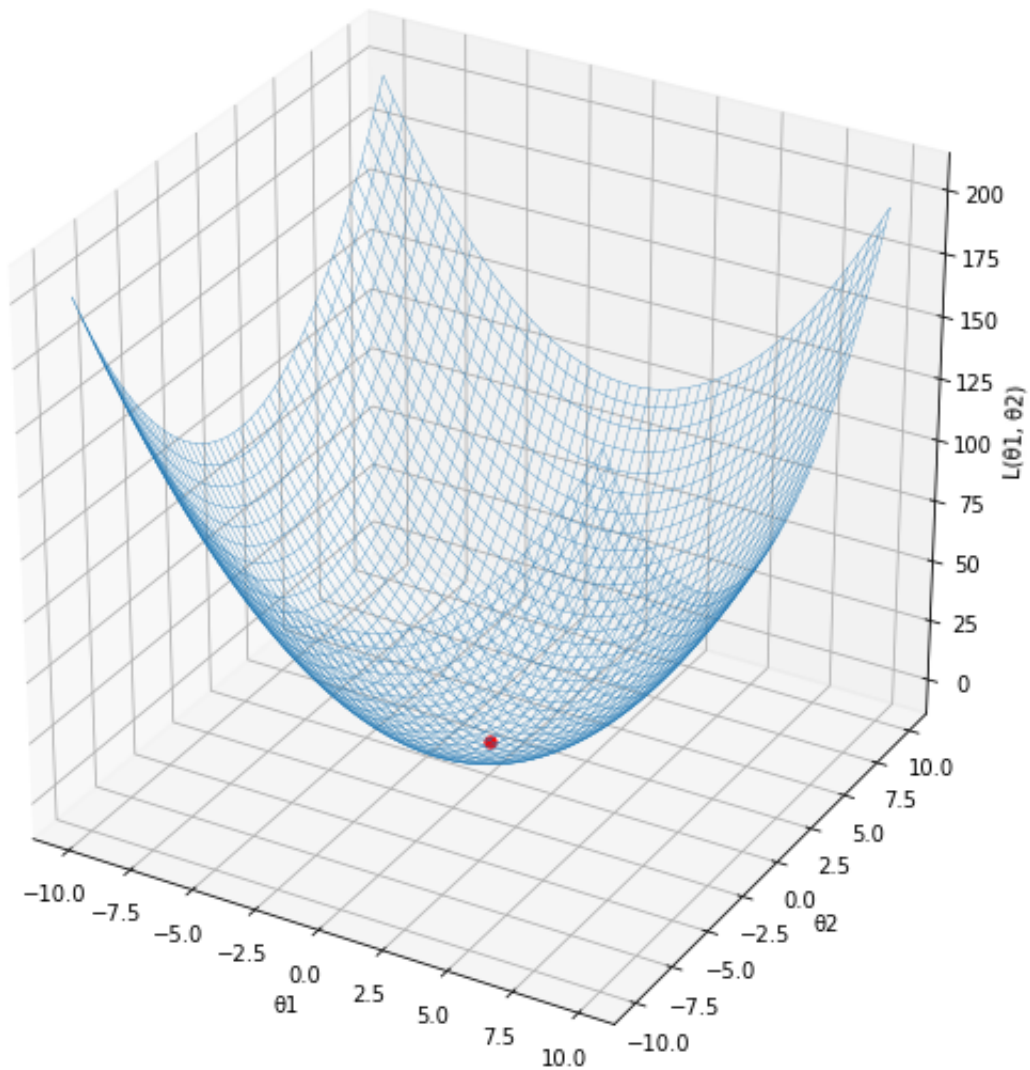
```

28 ax.set_xlabel(' 1 ')
29 ax.set_ylabel(' 2 ')
30 ax.set_zlabel('L( 1 , 2 )')
31 plt.show()

```

Listing 3: Question 2

Result



Observation/ Justification

The minimum value of $L(\theta)$ is — at —

Question 3(a)

Plot for $L(\theta) = \sum_{i=1}^m (y(i) - (\theta_0 + \theta_1 \cdot x(i)))^2$, where m is the number of input examples and $x(i)$, $y(i)$ are the values taken from the given data file. Obtain the minimum value of $L(\theta)$ with corresponding θ_0, θ_1 values from the plot.

Answer

code

```

1 def L(theta):
2     return theta**2
3
4 theta = np.arange(-10, 10, 0.1)
5 func_L = L(theta)
6 plt.figure()
7 plt.plot(theta, func_L, linewidth=2)
8
9 min_func_L = min(func_L)
10 min_location = np.argwhere(func_L == min_func_L)
11 #print(min_location)
12 plt.plot(theta[min_location[0,0]], min_func_L, 'ro', markersize=10)
13 plt.xlabel(' ')
14 plt.ylabel('L( )')
15 plt.show()

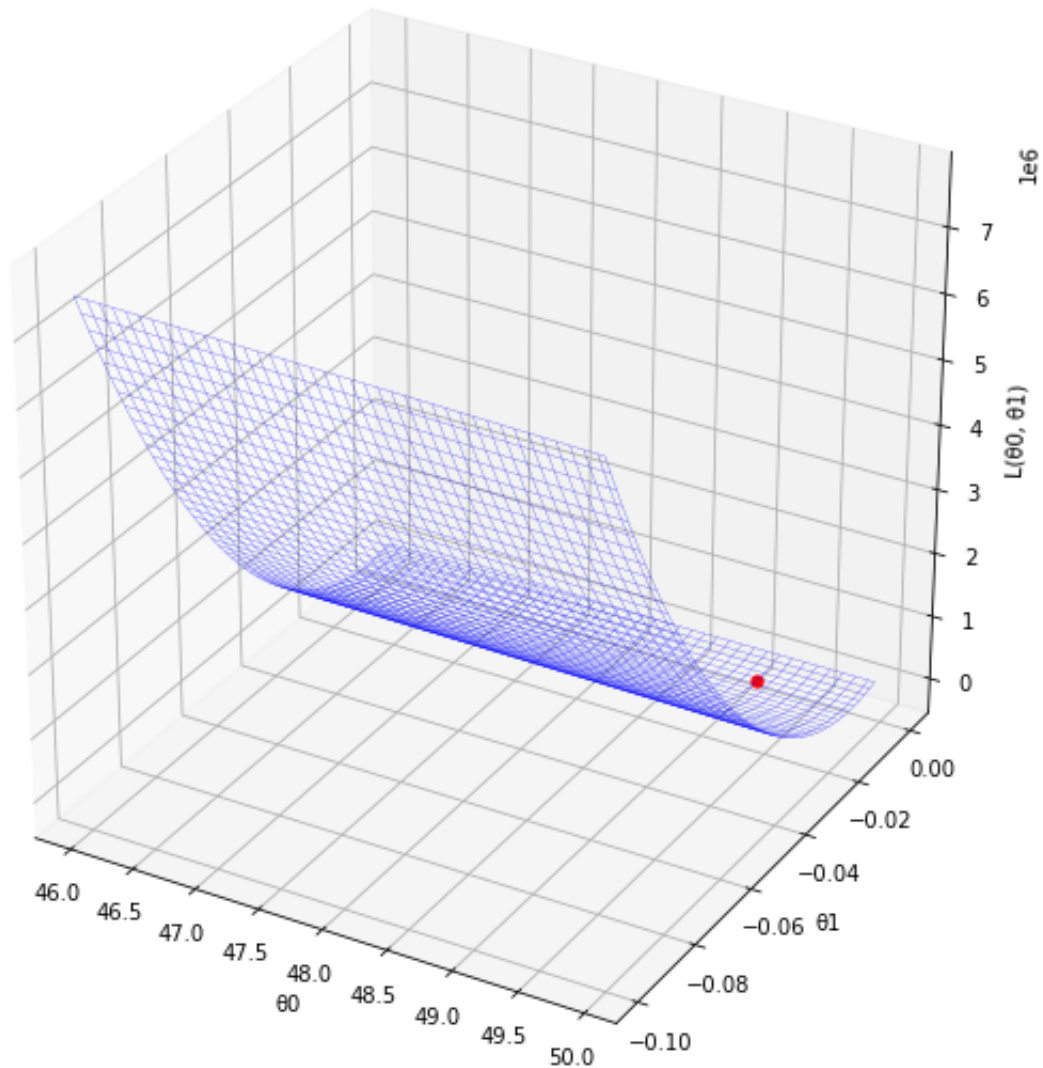
```

Listing 4: Question 3(a)

Result

OUTPUT:

Minimum value of $L(\theta_0, \theta_1)$ is 1572.6593617600001 θ_0 and θ_1 values for minimum value of $L(\theta_0, \theta_1)$ are 49.19999999999982 -0.008599999999997388 respectively.



Observation/ Justification

Question 3(b)

Apply Pseudo Inverse (Least Squares (LS)) approach to get θ vector for the cost function (objective function) $L(\theta)$ given in example 3(a). Verify whether θ_0, θ_1 obtained are the same as that found in Q3(a).

Answer

code

```
1 temp = np.copy(data_x)
2 x = np.ones((np.size(temp,0),2))
3 x[:,1] = temp[:,0]
4
5 print('determinant of x.T*x = ',np.linalg.det(np.matmul(x.T,x)))
6
7 theta = np.matmul(np.linalg.inv(np.matmul(x.T,x)),np.matmul(x.T,data_y))
8 print('For minima, theta0 = ',theta[0,0], ' and theta1 = ',theta[1,0])
```

Listing 5: Question 3(b)

Result

OUTPUT:

determinant of $x.T*x = 5618077959.999997$ For minima, $\theta_0 = 49.23762989433493$ and $\theta_1 = -0.008611934783475328$

Observation/ Justification

Question 4

Calculate the value of $L(\theta)$ using the θ vector obtained by Pseudo Inverse (as done in Q3(b)). Now Assume any θ vector (other than the one obtained in Q3(b)) and compute the new $L(\theta)$ value. Comment on why the Pseudo Inverse is also called LS method.

Answer

code

```
1 LS = np.sum(np.square(data_y - (theta[0,0] + (theta[1,0]*temp))))
2 print('Value of L using theta values from LS method(Pseudo Inverse) = ',LS)
3 print()
4 for i in range(10):
5     L_other_value = np.sum(np.square(data_y - (np.random.randint(theta0.shape[0]) - (np.random
6         .randint(theta1.shape[0])*temp))))
7     print('Value of L for some other theta = ',L_other_value)
```

Listing 6: Question 3(b)

Result

OUTPUT:

Value of L using theta values from LS method(Pseudo Inverse) = 1572.6503668922921

Value of L for some other theta = 2484947822720726.0 Value of L for some other theta = 269611008855161.75

Value of L for some other theta = 1140041054488897.8 Value of L for some other theta = 259697580389616.75

Value of L for some other theta = 1422578963889964.8 Value of L for some other theta = 12106098743706.75

Value of L for some other theta = 1385272296102789.8 Value of L for some other theta = 2036913111196393.8

Value of L for some other theta = 2838967870749331.0 Value of L for some other theta = 1151197388894299.8

Observation/ Justification

It can be observed that the value $L(\theta)$ for θ_0 and θ_1 found in the 4th Question gives a lesser value than other values found. Thus, it can be said that the Pseudo Inverse method gives the least values for the given function. Hence this is the reason why Pseudo Inverse is also called the Least Squares method.

Question 5(a)

For the following X and Y, use scikit-learn to learn a linear model.

$$X = \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \\ 7 & 8 \end{bmatrix} \quad Y = \begin{bmatrix} 2 \\ 3 \\ 4 \\ 5 \end{bmatrix}$$

Answer

code

```
1 from sklearn.linear_model import LinearRegression
2 from sklearn.metrics import mean_squared_error
3
4 x = np.array([
5     [1, 2],
6     [2, 4],
7     [3, 6],
8     [4, 8]
9 ])
10 y = np.array([
11     [2],
12     [3],
13     [4],
14     [5]
15 ])
16
17 reg = LinearRegression().fit(x,y)
18 print('Score: ', reg.score(x, y))
19 print('Coefficient: ', reg.coef_)
20 print('Intercept: ', reg.intercept_)
```

Listing 7: Question 5(a)

Result

Score: 1.0 Coefficient: [[0.2 0.4]] Intercept: [1.]

Observation/ Justification

Question 5(b)

Solve the problem using normal equations. You may find that one of the matrices in the normal equation is non-invertible. Why does the matrix turn out to be non-invertible? Why can scikit-learn implementation still correctly solve this regression problem?

Answer

Observation/ Justification

We got Singular Matrix Error because $\det(X'X) = 0$. Even though, scikit learn can solve the model as it uses pseudo inverse instead of normal inverse. So, if we use pinv fuction instead of inv function, we can still get the correct results.

Question 6(a)

Show the usage of scikit-learn's linear regression module for the real estate price prediction regression problem. What is the RMS error on the test set?

Answer

code

```

1 real_estate_data = pd.read_excel('/content/sample_data/real_estate_valuation_dataset.xlsx')
2 #real_estate_data
3 x1 = np.reshape(np.array(real_estate_data['X1 transaction date']), (-1, 1))
4 x2 = np.reshape(np.array(real_estate_data['X2 house age']), (-1, 1))
5 x3 = np.reshape(np.array(real_estate_data['X3 distance to the nearest MRT station']), (-1, 1))
6 x4 = np.reshape(np.array(real_estate_data['X4 number of convenience stores']), (-1, 1))
7 x5 = np.reshape(np.array(real_estate_data['X5 latitude']), (-1, 1))
8 x6 = np.reshape(np.array(real_estate_data['X6 longitude']), (-1, 1))
9 X = np.hstack((x1, x2, x3, x4, x5, x6))
10 y = np.reshape(np.array(real_estate_data['Y house price of unit area']), (-1, 1))
11
12 reg = LinearRegression().fit(X, y)
13 print('Score: ', reg.score(X, y))
14 print('Coefficient: ', reg.coef_)
15 print('Intercept: ', reg.intercept_)
16
17 y_pred1 = reg.predict(X)
18
19 print('RMSE ->', np.sqrt(mean_squared_error(y, y_pred1)))

```

Listing 8: Question 6(a)

Result

OUTPUT:

Score: 0.5823850447850936 Coefficient: [[5.14901721e+00 -2.69696735e-01 -4.48750825e-03 1.13332498e+00 2.25470143e+02 -1.24290612e+01]] Intercept: [-14441.98271918] RMSE -> 8.78231297536097

Observation/ Justification

Question 6(b)

Based on the regression coefficients, what can you comment about the importance of different features? Is it correct to assume that larger coefficients mean more important features?

Answer

Observation/ Justification

We know that the sign of a regression coefficient indicates whether there is a positive or negative correlation between each dependent variable and independent variable. A positive coefficient suggests that as the value of the independent variable increases, the mean of the dependent variable also tends to increase. A negative coefficient indicates that as the independent variable increases, the dependent variable tends to decrease.

The coefficient value signifies how much the mean of the dependent variable changes when the independent variable is shifted by one-unit while the other variables in the model are kept constant. By keeping the other variables constant it allows us to assess the effect of each variable in isolation from the others.

Although since the X values are not standardized, we cannot come to a conclusion on the importance of a feature by just comparing the coefficients.

Question 6(c)

Now, standardize the dataset to have all features on a scale of 0 to 1. Re-learn the regression coefficients and now comment on the importance of different features.

Answer

code

```

1 X_norm = (X - X.min())/(X.max()-X.min())
2 reg = LinearRegression().fit(X_norm, y)
3 print('Score: ', reg.score(X_norm, y))
4 print('Coefficient: ', reg.coef_)
5 print('Intercept: ', reg.intercept_)
6
7 y_pred2 = reg.predict(X_norm)
8
9 print('RMSE ->', np.sqrt(mean_squared_error(y, y_pred2)))

```

```

10
11 # creating bar graph
12 fig = plt.figure(figsize = (10, 5))
13 plt.bar([1,2,3,4,5,6], np.log10(abs(reg.coef_[0])), color = 'maroon',
14         width = 0.4)
15
16 plt.xlabel("Feature")
17 plt.ylabel("Score(log10)")
18 plt.title("Score for each feature")
19 plt.show()

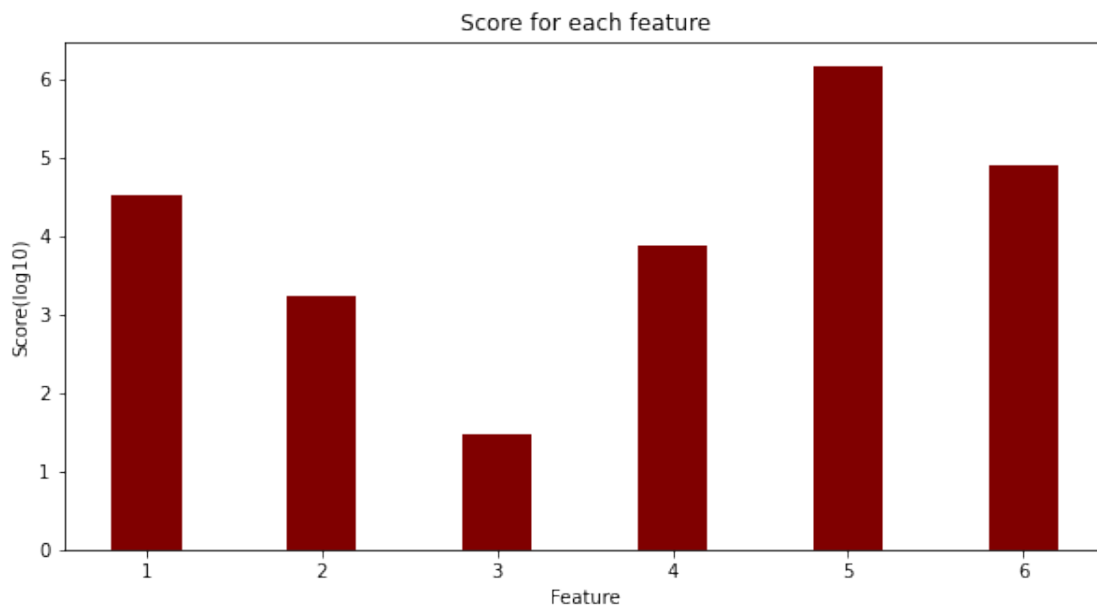
```

Listing 9: Question 3(b)

Result

OUTPUT:

Score: 0.5823850447850758 Coefficient: [[3.34069318e+04 -1.74979808e+03 -2.91150478e+01 7.35303628e+03
1.46285502e+06 -8.06400099e+04]] Intercept: [-14441.98271918] RMSE = 8.782312975361156



Observation/ Justification

Adding on to the observations made in 6(a), we can now say that after being normalized, the coefficients can be used as a measure of the importance of each feature w.r.t. the dependent variable. As seen in the graph, the 5th and 6th features are the most important ones, i.e., the latitude and longitude of the house matters the most. Moreover, the 3rd feature is the least important.

Question 6(e)

Answer

code

```

1 X = np.hstack((x1, x2, x3, x4, x5))
2 y = np.reshape(np.array(real_estate_data['Y house price of unit area']), (-1, 1))
3
4 reg = LinearRegression().fit(X, y)
5 print('Score: ', reg.score(X, y))
6 print('Coefficient: ', reg.coef_)
7 print('Intercept: ', reg.intercept_)
8
9 y_pred1 = reg.predict(X)
10
11 print('RMSE ->', np.sqrt(mean_squared_error(y, y_pred1)))

```

Listing 10: Question 6e)

Result

OUTPUT:

Score: 0.5823178829759132 Coefficient: [[5.13755503e+00 -2.69380517e-01 -4.35333829e-03 1.13619277e+00 2.26879404e+02]] Intercept: [-15964.80386271] RMSE : 8.783019142977878

Observation/ Justification

Upon removing the 6th feature, the error reduces. Thus the 6th feature is not that important. Similarly we can try for other features too and we will get different results for different features.