

# Programming Exercise 1: Linear Regression

Thiago Raulino Dal Pont  
INE6116000 - Conexión Artificial Intelligence  
Federal University of Santa Catarina

May 3, 2022

## 1 Introduction

In this first programming exercise, we are going to focus on Linear regression with one or more input variables.

Based on the Python Notebook available in Moodle, we implement some core functions for training Linear Regression models based on Gradient Descent and the Normal Equation method.

This work is divided as follows:

- Section 2 presents the required materials, methods and studies to finish the exercise;
- Section 3 presents the implementations, results and discussions;
- Section 4 presents the conclusions.

## 2 Methods and Materials

Before diving in the code, the exercise required some previous study in the Python libraries: `numpy` and `matplotlib`. To fulfill this goal, the documentations [1, 2] and some Youtube videos [3, 4] were used. A study on the Linear Regression using the reference book [5] was also required.

Although model evaluation is not required by the exercise, we applied two metrics to estimate how well the models fit the train data:

- Root Mean Square Error (RMSE):

$$RMSE = \sqrt{\frac{\sum_{i=1}^n (y_{act} - y_{pred})^2}{n}} \quad (1)$$

- $R^2$  - Coefficient of Determination

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_{act} - y_{pred})^2}{\sum_{i=1}^n (y_{act} - \bar{y})^2} \quad (2)$$

In terms of tools versions, we used:

- Python 3.9.7
- Numpy 1.20.3
- Matplotlib 3.5.1

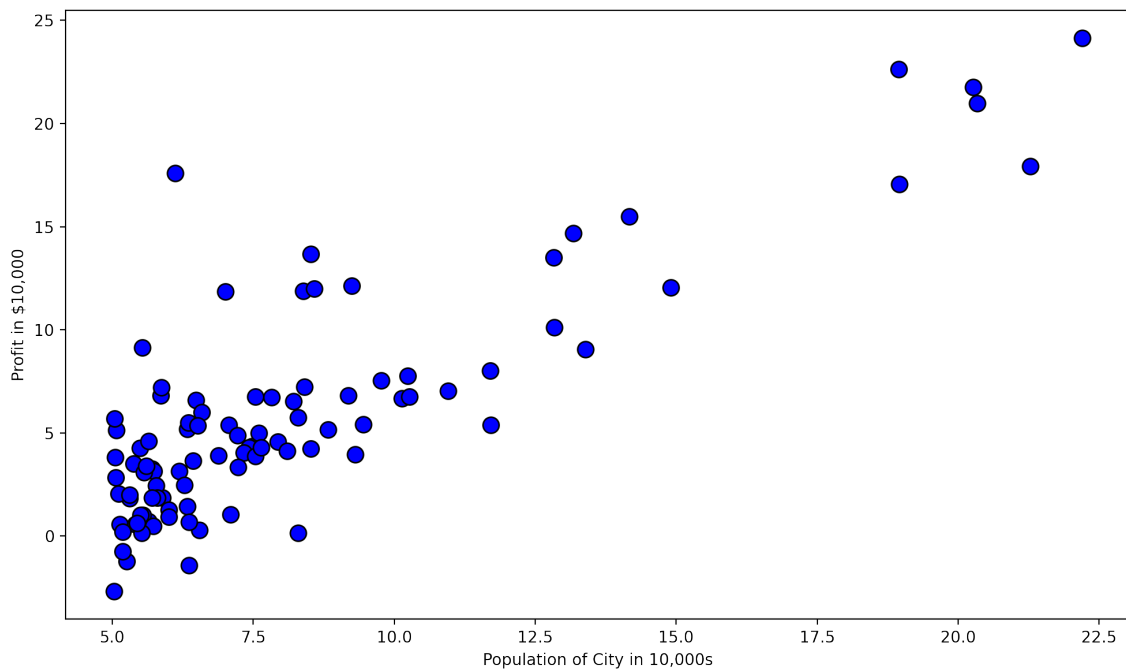
## 3 Results and Discussion

### 3.1 Linear regression with one variable

In this problem, we try to predict the profit of a food truck based on the city population.

The first step required in this section is the data visualization. Therefore, we implement the data plot using `pyplot` from `matplotlib` library, as shown in Figure 1.

Figure 1: Scatter plot of Population



The data shows a linear trend, which make its suitable for Linear Regression, according to Equation.

$$h_{\theta}(x) = \theta_0 + \theta_1 x \quad (3)$$

The next task was to compute the cost function  $J(\theta)$ . Using numpy functions, we could easily apply linear algebra operations on matrices.

The code is presented below.

Listing 1: Cost Function implementation

```
1 import numpy as np
2
3 def computeCost(X, y, theta):
4
5     # initialize some useful values
6     m = y.size # number of training examples
7
8     # You need to return the following variables correctly
9     J = 0
10
11     # ===== YOUR CODE HERE =====
12     ho_x = np.matmul(X, theta) # Calculate ho(x) for all elements in X
13     diff_mat = ho_x - y # Calculate the diff between all ho(x) and y
14     diff_mat_square = diff_mat ** 2 # Square the elements in the diff array
15     sum_op = np.sum(diff_mat_square) # Sum the elements of the squared array
16     J = (1 / (2 * m)) * sum_op # Multiply the sum with the constant
17     # =====
18
19     return J
```

Using the cost function, we implemented the Gradient Descent function

Listing 2: Cost Function implementation

```
1 def gradientDescent(X, y, theta, alpha, num_iters):
2     # Initialize some useful values
3     m = y.shape[0] # number of training examples
4
5     # make a copy of theta, to avoid changing the original array, since numpy arrays
6     # are passed by reference to functions
7     theta = theta.copy()
8
9     J_history = [] # Use a python list to save cost in every iteration
```

```

10
11 for i in range(num_iters):
12     # ===== YOUR CODE HERE =====
13     ho_x = X @ theta # [shape: (m,1)] Calculate ho(x) for all
14     elements in X
15     diff_mat = ho_x - y # [shape: (m,1)] Calculate the diff between all
16     ho(x) and y
17     mult_diff_xj = diff_mat @ X # [shape: (2,1)] Dot prod of diff with X matrix
18     const_mult = alpha * (1/m) # [shape: float] Get constant
19     theta = theta - const_mult * mult_diff_xj # [shape: [2,1]]: Update theta
20     values
21
22     # =====
23
24     # save the cost J in every iteration
25     J_history.append(computeCost(X, y, theta))
26
27 return theta, J_history

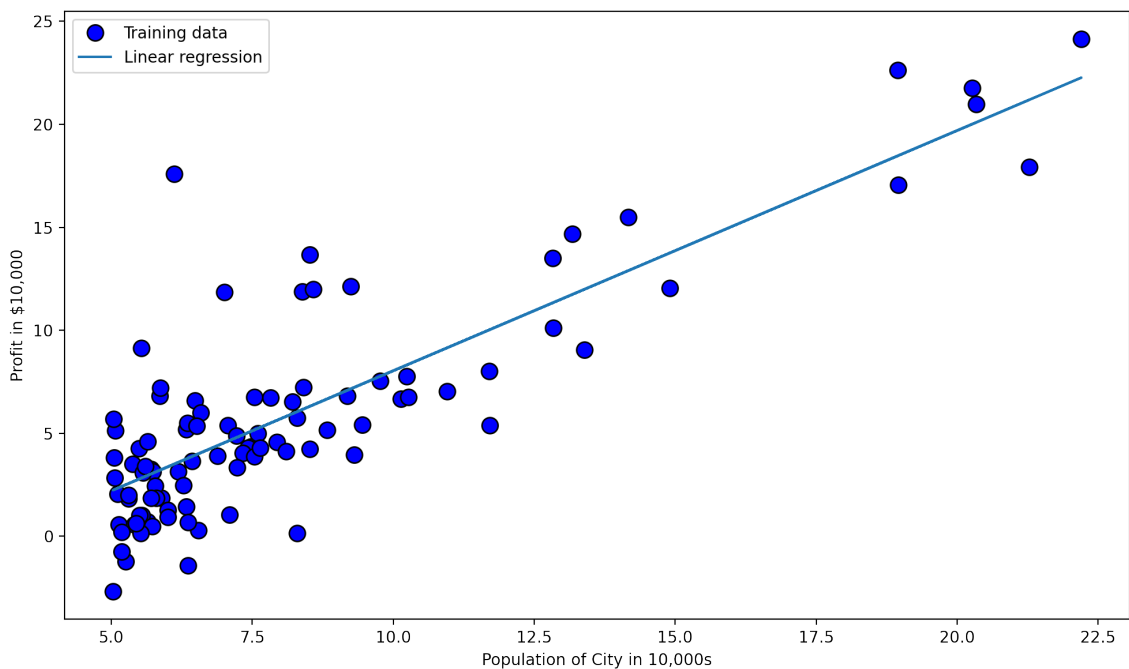
```

After training the GD using  $\alpha = 0.01$  with 1500 iterations, we obtained the following Equation:

$$h_{\theta}(x) = -3.6303 + 1.1664x \quad (4)$$

When plotting the hypothesis in Figure 2, we can see how well it fits the training data.

Figure 2: Plot hypothesis with training data



When calculating the regression metrics we obtained a RMSE of 2.99 and a  $R^2$  of 0.70.

One interesting exercise is to plot the cost function surface according to  $\theta_0$  and  $\theta_1$  from the Equation 3

### 3.2 Linear Regression with Multiple Variables

In the section, we try to predict house prices using the area ( $\text{ft}^2$ ) and the bedroom count.

Considering that the variables have distinct domain values, i.e., area is represented in thousands while bedroom in units, we need to normalize the dataset.

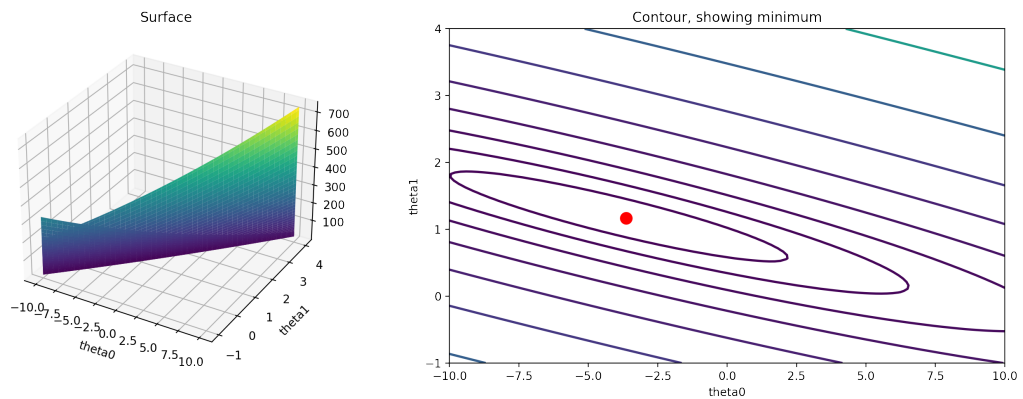
In the notebook, we applied the Standard Scaler, which represents the data as the z-score, using the formula:

$$X_{new} = \frac{X - \mu}{\sigma} \quad (5)$$

where  $\mu$  is the mean of  $X$  and  $\sigma$  the standard deviation.

The code that implements the normalization is presented in Listing 3.

Figure 3: Cost function surface



Listing 3: Feature Normalization

```
1 import numpy as np
2 def featureNormalize(X):
3
4     X_norm = X.copy()
5     # ===== YOUR CODE HERE =====
6     mu = np.mean(X_norm, axis=0)
7     sigma = np.std(X_norm, axis=0)
8
9     for i in range(len(mu)):
10         X_norm[:, i] = X_norm[:, i] - mu[i]
11         X_norm[:, i] = X_norm[:, i] / sigma[i]
12     # =====
13     return X_norm, mu, sigma
```

After implementing normalization, we adapted the cost function to allow multiple variable regression.

Listing 4: Cost function for LR with multiple variables

```
1 def computeCostMulti(X, y, theta):
2
3     # Initialize some useful values
4     m = y.shape[0] # number of training examples
5
6     # ===== YOUR CODE HERE =====
7     ho_x = np.matmul(X, theta) # Calculate ho(x) for all elements in X
8     diff_mat = ho_x - y # Calculate the diff between all ho(x)
9     # and y
10    diff_mat_square = np.dot(diff_mat.T, diff_mat) # Square the elements in
11    # the diff array
12    sum_op = np.sum(diff_mat_square) # Sum the elements of the squared array
13    J = (1 / (2 * m)) * diff_mat_square # Multiply the sum with the constant
14    # =====
15    return J
```

Finally, the GD algorithm for multi variables is presented in Listing 5.

Listing 5: GD for multi variable linear regression

```
1 def gradientDescentMulti(X, y, theta, alpha, num_iters):
2
3     # Initialize some useful values
4     m = y.shape[0] # number of training examples
5
6     # make a copy of theta, which will be updated by gradient descent
7     theta = theta.copy()
8
9     J_history = []
10
11    for i in range(num_iters):
12        # ===== YOUR CODE HERE =====
13        ho_x = X @ theta # [shape: (m,1)] Calculate ho(x) for all
14        # elements in X
15        diff_mat = ho_x - y # [shape: (m,1)] Calculate the diff between all
16        # ho(x) and y
```

```

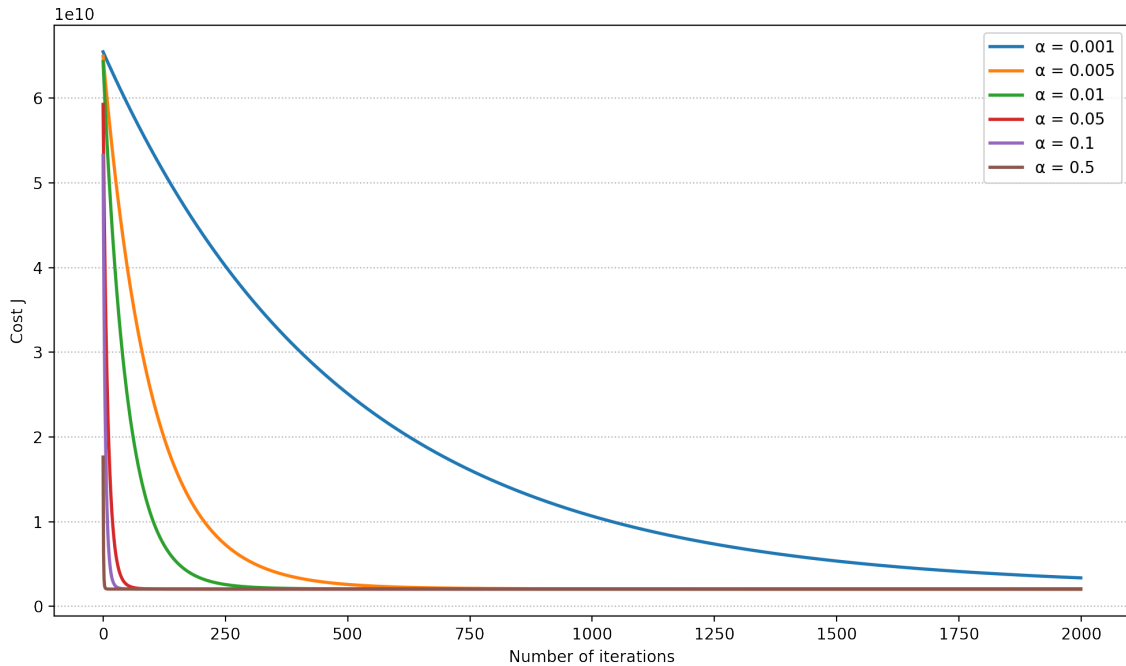
15     mult_diff_xj = diff_mat @ X    # [shape: (2,1)] Dot prod of diff with X matrix
16     const_mult = alpha * (1/m)    # [shape: float] Get constant
17     theta = theta - const_mult * mult_diff_xj # [shape: [2,1]]: Update theta
18     values
19     # =====
20     #print(X.shape, y.shape, ho_x.shape, diff_mat.shape, mult_diff_xj.shape,
21     theta.shape)
22     #print(theta)
23
24     J = computeCostMulti(X, y, theta)
25     # save the cost J in every iteration
26     J_history.append(J)
27
28     #if i % 100 == 0:
29     #    print("Passo: %d;\tJ=%.2f;\ttheta=%s" %(i, J, str(theta)))
30
31     #time.sleep(1)
32
33     return theta, J_history

```

Although the algorithm may work with several learning rates ( $\alpha$ ), we may improve the results as we search for the optimal value.

We tested several values for learning rate, as shown in Figure 4.

Figure 4: GD using several  $\alpha$  values



We may see GD may need from around 100 to more than 2000 iterations to converge depending on the learning rate.

If we consider both the number of iterations required for the convergence and the cost, we may estimate that a value of  $\alpha = 0.05$  is a good choice.

Table 1 presents the metrics after 2000 iterations for some values of learning rates.

### 3.3 Normal Equations

The values of  $\theta_0$  and  $\theta_1$  may be acquired analytically, using the following equation.

$$\theta = (X^T X)^{-1} X^T y \quad (6)$$

Using numpy library, such equation can be easily implemented as shown in Listing 6.

Table 1: Regression metrics for some values of learning rate

$\alpha$	RMSE	$R^2$
0.001	82037	0.71
0.005	63926	0.73
0.010	63926	0.73
0.050	63926	0.73
0.100	63926	0.73
0.500	63926	0.73

Listing 6: Normal Equation for multi variable linear regression

```

1 def normalEqn(X, y):
2     theta = np.zeros(X.shape[1])
3
4     # ===== YOUR CODE HERE =====
5     x_t = X.T
6     xt_x = np.matmul(x_t, X)
7     xt_x_inv = np.linalg.inv(xt_x)
8     xt_x_inv_xt = np.matmul(xt_x_inv, x_t)
9     theta = np.matmul(xt_x_inv_xt, y)
10    # =====
11    return theta

```

Using the normal equation, we achieved a RMSE of 63926, a  $R^2$  of 0.73. It is worth reminding that those are the same metrics achieved by the Gradient Descent algorithm, although the values of  $\theta^1$  are quite different.

## 4 Final Remarks

In this exercise, we focused on implementing the Linear Regression using Gradient Descent and the Normal Equation, with one and multiple variables.

Using the `numpy` library, applying linear algebra operations with matrices was straight forward.

In terms of values of  $\alpha$ , we saw that a low value makes it difficult for the algorithms to converge, while higher values may reduce the number of iterations required for convergence. However, although not presented in this exercise, even higher values of  $\alpha$  may blow up the the cost function values.

## References

- [1] "Mathematical functions - numpy." <https://numpy.org/doc/stable/reference/routines.math.html>. Accessed 30 Apr 2022.
- [2] "Matplotlib 3.5.1 documentation." <https://matplotlib.org/3.5.1/index.html>. Accessed 30 Apr 2022.
- [3] "Complete python numpy tutorial (creating arrays, indexing, math, statistics, reshaping)." <https://www.youtube.com/watch?v=GB9ByFAIAH4>. Accessed 30 Apr 2022.
- [4] "Intro to data visualization in python with matplotlib! (line graph, bar chart, title, labels, size)." <https://www.youtube.com/watch?v=DAQNHZocO5A>. Accessed 30 Apr 2022.
- [5] S. Haykin, *Neural networks and Learning Machines*. 2008.

---

<sup>1</sup>See details in the notebook