

---

# Exercício de Programação 1

## Regressão Linear

---

Sheldon Goulart de Alcântara  
alcantarash92@gmail.com  
Vinícius Teodoro  
vtcpires@gmail.com

### Abstract

This report aims to explain the linear regression and the equations used to obtain the results proposed in the exercise.

## 1 Regressão Linear

Modelos de regressão são modelos matemáticos que relacionam o comportamento de uma variável Y com uma variável X. Quando a função  $f$  que relaciona duas variáveis é do tipo  $f(X) = a + bX$  temos o modelo de regressão simples. A variável X é a variável independente da equação, enquanto  $Y = f(X)$  é dependente das variações de X. O modelo de regressão é chamado de simples quando envolve uma relação casual entre duas variáveis. O modelo de regressão é múltiplo quando envolve uma relação causal com mais de duas variáveis. Isto é, quando o comportamento de Y é explicado por mais de uma variável independente  $X_1, X_2, \dots, X_n$ .

Os modelos acima (simples ou multivariados) simulam relacionamentos entre as variáveis. Esse relacionamento poderá ser do tipo linear (equação da reta ou do plano) ou não linear (equação exponencial, geométrica, etc.). A análise de regressão compreende, portanto quatro tipos básicos de modelos:

- Linear simples.
- Linear multivariado.
- Não linear simples.
- Não linear multivariado.

### 1.1 Modelos de regressão linear

Regressão é o processo matemático pelo qual derivamos os parâmetros 'a' e 'b' de uma função  $f(X)$ . Estes parâmetros determinam as características da função que relaciona 'Y' com 'X' que no caso do modelo linear se representa por uma reta chamada de reta de regressão. Esta reta explica de forma geral e teoricamente a relação entre X e Y. Isto significa que os valores observados de X e Y nem sempre serão iguais aos valores de X' e Y' estimados pela reta de regressão. Haverá sempre alguma diferença, e essa diferença significa:

- As variações de Y não são perfeitamente explicadas pelas variações de X.
- Existem outras variáveis das quais Y depende.
- Os valores de X e Y são obtidos de uma amostra específica que apresenta distorções em relação a realidade.

Esta diferença em estatística é chamada de erro ou desvio.

O processo de regressão significa, portanto, traçar uma reta com menor distâncias entre os pontos plotados no gráfico, ou seja, diminuir a diferença entre  $Y$  e  $Y'$ .

$$y_0 = a_0 + b_0.x_0$$

equação da reta a partir dados coletados

$$y_1 = +b_1.x_1$$

equação da reta a partir das estimativas

## 2 Função de custo

A função de custo (ou função de erros) é a função que define a qualidade da linha de regressão. Essa função recebe um par ordenado  $(\beta_0, \beta_1)$  e retorna um valor de erro baseado no quão bem a reta se adequa aos dados a serem analisados. Ao utilizar o método da regressão linear temos por objetivo minimizar a função de custo:

$$J(\beta) = \frac{1}{2m} \sum_{i=1}^n (h_{\beta}(x_i) - y_i)^2 \quad (1)$$

Onde o modelo  $h_{\beta}$  é dado pela reta:

$$h_{\beta}(x) = \beta^T x = \beta_0 + \beta_1 x_1 \quad (2)$$

Na implementação em python utilizamos de operações matriciais da biblioteca NumPy para evitar a criação de loops desnecessários. A implementação da função de erros pode ser vista no Algoritmo 1 (compute cost):

```
1 def compute_cost(X, y, beta):
2     lengthX = 2 * len(X)
3     funct = np.power(((X * beta.T) - y), 2)
4     return np.sum(funct) / lengthX
```

Listing 1: Algoritmo 1 ComputeCost

## 3 Método Gradiente

O método do gradiente é usado em otimização. Para encontrar um mínimo (local) de uma função usa-se um esquema iterativo, onde em cada passo se toma a direção (negativa) do gradiente, que corresponde à direção de declive máximo.

### 3.1 Descrição

Começando com um vetor inicial  $x_0$  visando alcançar um ponto mínimo de  $f$ , consideramos a sucessão definida por:

$$x_0, x_1, x_2, \dots, x_n$$

onde a pesquisa linear é dada pela direção de descida  $d_n$

$$x_{n+1} = x_n + \omega_n.d_n$$

No caso do método do gradiente a condição de descida verifica-se tomando

$$d_n = -\nabla F(x_n)$$

ficando a iteração definida por

$$x_{n+1} = x_n - \omega \nabla F(x_n)$$

O algoritmo utilizado para a implementação da função Gradiente Descendente é relativamente simples. Dado um ponto inicial da superfície da função que se deseja minimizar ele encontra o sentido

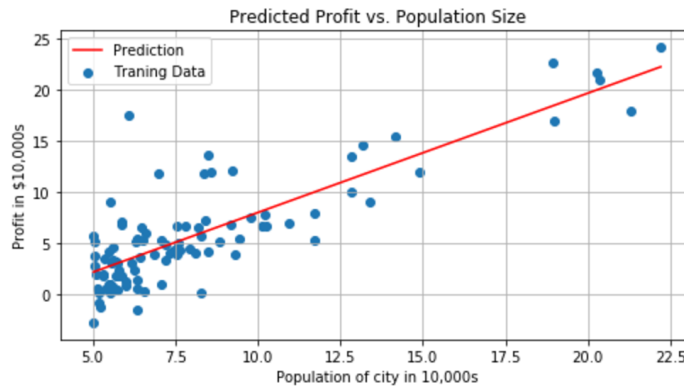


Figure 1: Função  $J(\beta)$  representada graficamente

de crescimento da função e se move em sua direção, segundo um dado parâmetro alfa ( $\alpha$ ), até atingir um ponto de convergência.

$$\beta_j = \beta_j - \alpha \frac{1}{n} \sum_{i=1}^n (h_{\beta}(x_i) - y_i) x_{ij} \quad (3)$$

```

1 def gradient_descent(X, y, theta, alpha, iters):
2     temp = np.matrix(np.zeros(theta.shape)) #np.zeros: Return a new
3     array of given shape and type, filled with zeros.
4     parameters = int(theta.ravel().shape[1]) #Return a contiguous
5     flattened array and return a tuple of array dimensions axis 1
6     cost = np.zeros(iters)
7     n = len(X)
8     for i in range(iters):
9         for j in range(parameters):
10             gradient = np.multiply((X.dot(theta.T) - y), X[:, j])
11             temp[0, j] -= (alpha / (n)) * np.sum(gradient)
12
13         theta = temp #update
14         cost[i] = compute_cost(X, y, theta)
15     return theta, cost

```

Listing 2: Algoritmo 2 GradientDescent

Pode-se observar no algoritmo 2 (GradientDescent), a implementação em python dessa função. Note que as variáveis temp0 e temp1 são necessárias para que os parâmetros ( $\beta_0$ ,  $\beta_1$ ) sejam atualizados de forma simultânea. De modo contrário o cálculo dos parâmetros podem se influenciar mutuamente, gerando um resultado incorreto.

## 4 Regressão Linear Múltipla

As funções desenvolvidas nesse trabalho podem ser usadas tanto para uma regressão linear simples (uma variável), como para regressões lineares múltiplas. Para demonstrar isso utilizaremos o DataSet sobre preços de imóveis como exemplo.

Inicialmente pode-se observar a relação entre o número de quartos e o tamanho do imóvel com seu preço plotando os dados do dataset em um gráfico de dispersão. Porém, antes de mais nada, note que como a dimensão das variáveis é bem diferente. Comparando os tamanhos das casas com seu número de quartos e o preço, pode-se ver como a discrepância entre esses números é grande. Para solucionar esse problema utilizou-se um método chamado 'feature normalization'. Ele consiste em subtrair de cada atributo, o valor da média de sua coluna, e dividi-lo pelo seu desvio padrão. Desta forma podemos ajustar as dimensões utilizadas.

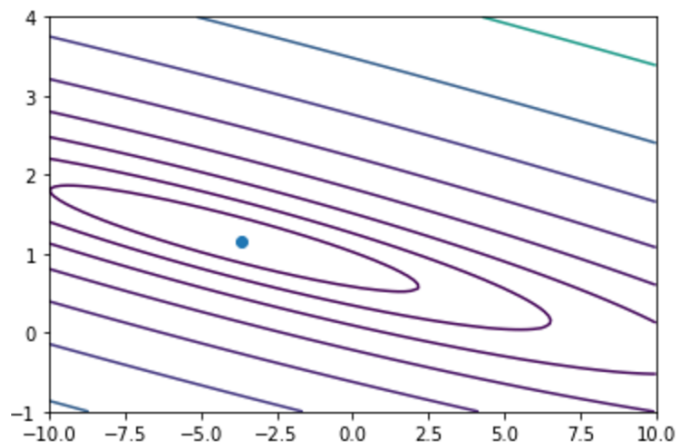


Figure 2: Gráfico do dataset com as dimensões devidamente ajustadas

Após isso podemos comparar o número de iterações de nosso algoritmo com os valores retornados de nossa função de custos

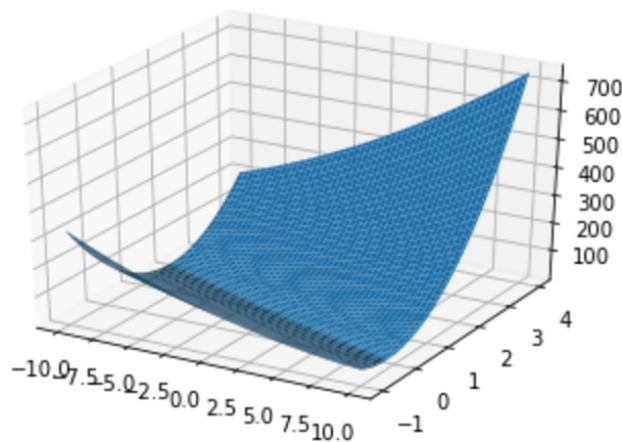


Figure 3: Número de iterações de nosso algoritmo com os valores retornados da função da custos

Pode-se visualizar todo o código de estudo utilizado em algoritmo 3 (Code Data2) e um exemplo do mundo real no algoritmo 4 (Code Groupon) :

```

1 # Load dataset
2 data2 = pd.read_csv('ex1data2.txt', header=None, names=['Size', 'Bedrooms',
3               'Price'])
4
5 # Executing the 'feature normalization'
6 data2 = (data2 - data2.mean()) / data2.std()
7 data2.head()
8
9 data2.insert(0, 'beta zero', 1)
10
11 # set X (training data) and y (target variable)
12 cols = data2.shape[1]
13 x = data2.iloc[:, 0:cols]
```

```

216 13 y = data2.iloc[:, cols-1:cols]
217 14
218 15 x = np.matrix(x.values)
219 16 y = np.matrix(y.values)
220 17 beta2 = np.matrix(np.array([0,0,0]))
221 18 g2, cost2 = gradient_descent(x, y, beta2, alpha, iters)
222 19
223 20 x = np.linspace(data2.Size.min(), data2.Size.max(), 100)
224 21 f = g2[0, 0] + (g2[0, 1] * x) + (g2[0, 2] * x)
225 22
226 23 fig, ax = plt.subplots(figsize=(8,4))
227 24 ax.plot(x, f, 'r', label='Prediction')
228 25 ax.scatter(data2.Size, data2.Price, label='Traning Data2')
229 26 ax.legend(loc=2)
230 27 ax.set_xlabel('Size')
231 28 ax.set_ylabel('Price')
232 29 ax.set_title('Predicted Price vs. Size')
233 30 ax.grid(True)
234 31
235 32 fig, ax = plt.subplots(figsize=(12,8))
236 33 ax.plot(np.arange(iters), cost2, 'r')
237 34 ax.set_xlabel('Iterations')
238 35 ax.set_title('Error vs. Training Epoch')
239 36 ax.grid(True)
240 37
241 38
242 39 beta0_vals = np.linspace(-10, 10, 100)
243 40 beta1_vals = np.linspace(-1, 4, 100)
244 41
245 42 for i in range(len(beta0_vals)):
246 43     for j in range(len(beta1_vals)):
247 44         t = np.matrix(np.array([beta0_vals[i], beta1_vals[j]]))
248 45         j_vals[i,j] = compute_cost(X, y, t)
249 46
250 47
251 48 plt.scatter(g2[0,0], g2[0,1],)
252 49 plt.contour(beta0_vals, beta1_vals, j_vals.T, np.logspace(-2, 3, 20));
253 50
254 51 beta0_mesh, beta1_mesh = np.meshgrid(beta0_vals, beta1_vals)
255 52 fig = plt.figure()
256 53 ax = fig.gca(projection='3d')
257 54 ax.plot_surface(beta0_mesh, beta1_mesh, j_vals.T);

```

Listing 3: Code Data2

```

254 1 # Load dataset
255 2 data3 = pd.read_csv('groupon-deals.csv', header=0)
256 3
257 4 # Make qualitative variables quantitative
258 5 data3['city'] = pd.Categorical(data3.city).codes
259 6 data3['category'] = pd.Categorical(data3.category).codes
260 7 data3['weekday_start'] = pd.Categorical(data3.weekday_start).codes
261 8
262 9
263 10 x = pd.concat([data3.iloc[:,1:3], data3.iloc[:,11:15]], axis=1)
264 11 x.insert(0, 'beta zero', 1)
265 12 y = data3.iloc[:,5:6]
266 13
267 14 x = np.matrix(x.values)
268 15 y = np.matrix(y.values)
269 16 theta = np.zeros(shape=(1,7))
270 17
271 18 # Plot cost function results
272 19 alpha = 1.00
273 20 iters = 100

```

```

270
271 21 theta , J_history = gradient_descent(x, y, theta , alpha , iters)
272 22 plt.plot(np.arange(iters), J_history)
273 23 plt.xlabel('Iterations')
274 24 plt.ylabel('Cost Function')
275 25 plt.show()
276 26

```

Listing 4: Code Groupon

```

276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

```