

Entrega 2: regresión lineal de multivariable

Aprendizaje Automatico y Big Data- Alejandro Barrachina Argudo

Introducción

En este documento se explicará el código del entregable 2 y el proceso de la regresión lineal en variables múltiples.

Para esta práctica se usarán los siguientes *imports* vistos en la figura 0.1.

```
1 import numpy as np
2 import copy
3 import matplotlib.pyplot as plt
4 import public_tests
5 import utils
6 import csv
```

Figura 0.1: Código de las bibliotecas usadas

1. Regresión lineal en variables múltiples

Para esta regresión definimos la función de pendiente como en la figura 1.1 como métrica, también se usan las funciones de coste 1.2 y de gradiente 1.3 según descritas en los apuntes.

Para hacer el gradiente descendente usamos 1.4 con 1000 iteraciones, 0.01 en α y 0 para iniciar w y b .

Si no hacemos una normalización de los datos, estos procesos darán pesos bastante dispares e incorrectos a los distintos campos de predicción, para ello usamos la función 'zscore_normalize_features' (figura 1.5) que normaliza los datos de la entrada dados.

```
1 def fun(x: np.ndarray, w: np.ndarray, b: float) -> float:
2     """Returns a predicted w value for a linear regression with the given w and b values.
3
4     Args:
5         x (np.ndarray): x value for de function
6         w (np.ndarray): array of w
7         b (float): b value
8
9     Returns:
10         float: predicted y value
11     """
12
13     return np.dot(x, w) + b
```

Figura 1.1: Código de la función 'fun'

```

1 def compute_cost(X: np.ndarray, y: np.ndarray, w: np.ndarray, b: float) -> float:
2     """
3     compute cost
4     Args:
5         X (ndarray (m,n)): Data, m examples with n features
6         y (ndarray (m,)) : target values
7         w (ndarray (n,)) : model parameters
8         b (scalar)       : model parameter
9     Returns
10        cost (scalar)    : cost
11    """
12    x_fun = fun(X, w, b)
13
14    cost = np.sum((x_fun - y)**2)
15
16    cost = cost / (2 * X.shape[0])
17
18    return cost

```

Figura 1.2: Código de la función 'compute_cost'

```

1 def compute_gradient(X: np.ndarray, y: np.ndarray, w: np.ndarray, b: float) -> tuple[float, np
2     .ndarray]:
3     """
4     Computes the gradient for linear regression
5     Args:
6         X : (ndarray Shape (m,n)) matrix of examples
7         y : (ndarray Shape (m,)) target value of each example
8         w : (ndarray Shape (n,)) parameters of the model
9         b : (scalar)             parameter of the model
10    Returns
11        dj_dw : (ndarray Shape (n,)) The gradient of the cost w.r.t. the parameters w.
12        dj_db : (scalar)             The gradient of the cost w.r.t. the parameter b.
13    """
14
15    x_fun = fun(X, w, b)
16    dj_dw = np.dot((x_fun - y), X) / X.shape[0]
17
18    dj_db: float = np.sum(x_fun - y)
19
20    dj_db /= X.shape[0]
21
22    return dj_db, dj_dw

```

Figura 1.3: Código de la función 'compute_gradient'

```

1 def zscore_normalize_features(X: np.ndarray) -> tuple[np.ndarray, np.ndarray, np.ndarray]:
2     """
3     computes X, zscore normalized by column
4
5     Args:
6         X (ndarray (m,n)) : input data, m examples, n features
7
8     Returns:
9         X_norm (ndarray (m,n)): input normalized by column
10        mu (ndarray (n,)) : mean of each feature
11        sigma (ndarray (n,)) : standard deviation of each feature
12    """
13    mu = np.mean(X, axis=0)
14    sigma = np.std(X, axis=0)
15    X_norm = (np.array(X) - mu) / sigma
16
17    return (X_norm, mu, sigma)

```

Figura 1.5: Código de la función ‘zscore_normalize_features’

```

1 def gradient_descent(X: np.ndarray, y: np.ndarray, w_in: np.ndarray, b_in: float,
2   cost_function: float,
3   gradient_function: float, alpha: float, num_iters: int) -> tuple[np.
4   ndarray, float, list[float]]:
5     """
6     Performs batch gradient descent to learn theta. Updates theta by taking
7     num_iters gradient steps with learning rate alpha
8
9     Args:
10        X : (array_like Shape (m,n)) matrix of examples
11        y : (array_like Shape (m,)) target value of each example
12        w_in : (array_like Shape (n,)) Initial values of parameters of the model
13        b_in : (scalar) Initial value of parameter of the model
14        cost_function: function to compute cost
15        gradient_function: function to compute the gradient
16        alpha : (float) Learning rate
17        num_iters : (int) number of iterations to run gradient descent
18    Returns
19        w : (array_like Shape (n,)) Updated values of parameters of the model
20        after running gradient descent
21        b : (scalar) Updated value of parameter of the model
22        after running gradient descent
23        J_history : (ndarray): Shape (num_iters,) J at each iteration,
24        primarily for graphing later
25    """
26
27    w = copy.deepcopy(w_in)
28    b = 0 + b_in
29    J_history = [cost_function(X, y, w, b)]
30
31    for i in range(0, num_iters):
32        (dj_db, dj_dw) = gradient_function(X, y, w, b)
33
34        w = w - (alpha * dj_dw)
35        b = b - (alpha * dj_db)
36        J_history.append(cost_function(X, y, w, b))
37
38    return w, b, J_history

```

Figura 1.4: Código de la función ‘gradient_descent’

2. Gráficas

La primera gráfica a generar es la que nos muestra los datos brutos (figura 2.5), viendo así como se relaciona cada característica con el precio de una vivienda, para ello usamos la función ‘visualize_data’ 2.1. Tras el entrenamiento del modelo podemos ver en la figura 2.6 la predicción de precios y los datos dados, este gráfico lo generamos con la función ‘visualize_data_train’ (figura 2.2).

También podemos observar la evolución de la función de coste en la figura 2.7 con la función ‘visualize_J_history’ (figura 2.4) y las regresiones parciales en la figura 2.8, producidas por la función ‘visualize_partial_regression’ (figura 2.3).

```

1 def visualize_data(X_train: np.ndarray, y_train: np.ndarray, name: str) -> None:
2     """Generates a scatter plot of the data
3
4     Args:
5         X_train (np.ndarray): training data
6         y_train (np.ndarray): training data
7         name (str): name of the plot
8     """
9     X_features = ['size(sqft)', 'bedrooms', 'floors', 'age']
10    fig, ax = plt.subplots(1, 4, figsize=(25, 5), sharey=True)
11    for i in range(len(ax)):
12        ax[i].scatter(X_train[:, i], y_train)
13        ax[i].set_xlabel(X_features[i])
14        ax[0].set_ylabel("Price (1000's)")
15    plt.savefig(f'./memoria/images/{name}.png', dpi=300)

```

Figura 2.1: Código de la función ‘visualize_data’

```

1 def visualize_data_train(X_train: np.ndarray, y_train: np.ndarray, prediction: np.ndarray) ->
None:
2     """Generates a scatter data with the original data and the predictions
3
4     Args:
5         X_train (np.ndarray): Train values
6         y_train (np.ndarray): Train prices
7         prediction (np.ndarray): Predicted prices
8     """
9     X_features = ['size(sqft)', 'bedrooms', 'floors', 'age']
10    fig, ax = plt.subplots(1, 4, figsize=(25, 5), sharey=True)
11    for i in range(len(ax)):
12        ax[i].scatter(X_train[:, i], y_train)
13        ax[i].scatter(X_train[:, i], prediction, color='orange')
14        ax[i].set_xlabel(X_features[i])
15        ax[0].set_ylabel("Price (1000's)")
16    plt.savefig(f'./memoria/images/predicted_data.png', dpi=300)

```

Figura 2.2: Código de la función ‘visualize_data_train’

```

1 def visualize_partial_regressions(X_train: np.ndarray, y_train: np.ndarray, w: np.ndarray, b:
float) -> None:
2     """Draws partial regressions for each variable
3
4     Args:
5         X_train (np.ndarray): Train data
6         y_train (np.ndarray): Train prices
7         w (np.ndarray): Weights
8         b (float): bias
9     """
10    (X_norm, _, _) = zscore_normalize_features(X_train)
11    X_features = ['size(sqft)', 'bedrooms', 'floors', 'age']
12    fig, ax = plt.subplots(1, 4, figsize=(25, 5), sharey=True)
13    for i in range(len(ax)):
14        ax[i].scatter(X_train[:, i], y_train)
15        ax[i].plot(X_train[:, i], fun(
16            X_norm[:, i], w[i], b), color='red')
17        ax[i].set_xlabel(X_features[i])
18        ax[0].set_ylabel("Price (1000's)")
19    plt.savefig(f'./memoria/images/partail_regression.png', dpi=300)

```

Figura 2.3: Código de la función ‘visualize_partial_regression’

```

1 def visualize_J_history(J_history: np.ndarray) -> None:
2     """Generates a plot with the evolution of the cost function
3
4     Args:
5         J_history (np.ndarray): Cost function over each iteration
6     """
7     plt.clf()
8     plt.figure(figsize=(7, 5))
9     plt.plot(J_history)
10    plt.xlabel('Iterations')
11    plt.xscale('log')
12    plt.ylabel('Cost')
13    plt.title('Cost function over iterations')
14    plt.savefig(f'./memoria/images/cost_function.png', dpi=300)

```

Figura 2.4: Código de la función 'visualize_J_history'

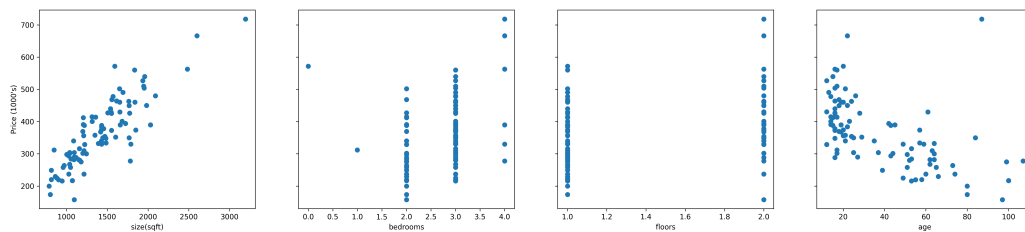


Figura 2.5: Visualización de los datos iniciales

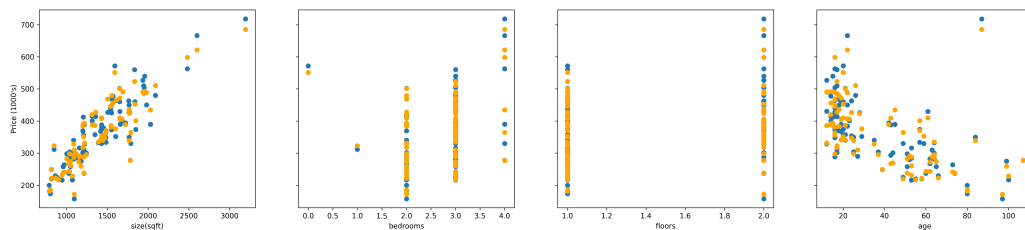


Figura 2.6: Visualización de los datos predichos

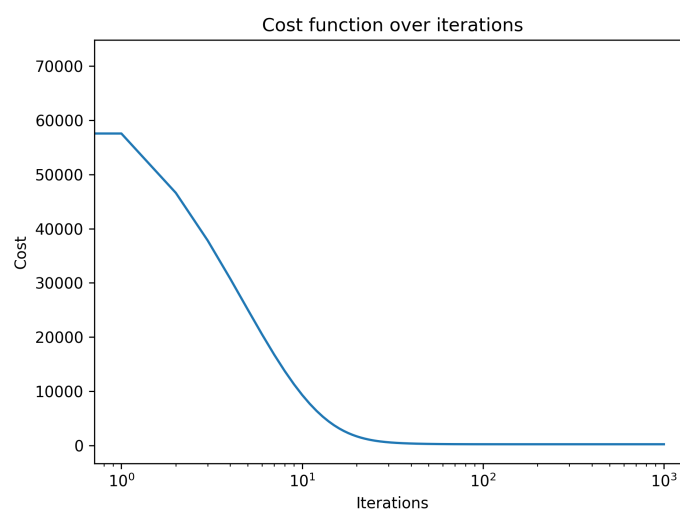


Figura 2.7: Evolución de la función de coste

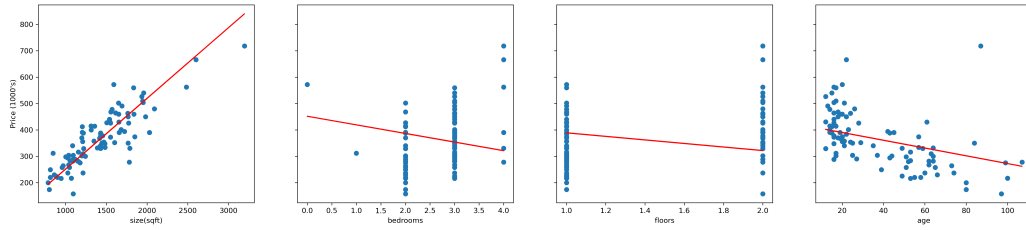


Figura 2.8: Regresiones parciales

3. conclusiones

Este modelo al tener más variables produce mejores predicciones que el modelo de la práctica anterior, ya que el anterior solo tenía en cuenta una variable arbitraria.

Para tener más presentes los resultados, usamos la función 'write_results' (figura 3.1) que nos escribe en un archivo csv los resultados de la predicción. Los resultados finales son:

| Iteraciones | $J(w,b)$ |
|-------------|--------------------|
| 100 | 71280.37305561663 |
| 200 | 221.21399549940065 |
| 300 | 219.20904207205783 |
| 400 | 219.2068268334026 |
| 500 | 219.20682438520015 |
| 600 | 219.20682438249446 |
| 700 | 219.20682438249148 |
| 800 | 219.2068243824914 |
| 900 | 219.2068243824914 |
| 1000 | 219.20682438249148 |

Cuadro 3.1: Evolución de $J(w,b)$

| size(sqft) | Bedrooms | floors | age |
|--------------------|-------------------|--------------------|-------------------|
| 1418.3737373737374 | 2.717171717171717 | 1.3838383838383839 | 38.38383838383838 |

Cuadro 3.2: Medias de cada atributo para la normalización

| size(sqft) | Bedrooms | floors | age |
|-------------------|--------------------|--------------------|--------------------|
| 411.6156289269652 | 0.6519652348729262 | 0.4863193178671001 | 25.777880687549874 |

Cuadro 3.3: Desviación estándar de cada atributo para la normalización

| W. size(sqft) | W. Bedrooms | W. floors | W. age | b |
|--------------------|---------------------|---------------------|--------------------|--------------------|
| 110.56039755974356 | -21.267150958179855 | -32.707181391704225 | -37.97015909100321 | 363.15608080808056 |

Cuadro 3.4: Pesos y bias finales para el modelo

```
1 def write_results(J_history: np.ndarray, w: np.ndarray, b: float, mu: np.ndarray, sigma: np.
  ndarray) -> None:
2   X_features = ['size(sqft)', 'bedrooms', 'floors', 'age']
3   with open('./memoria/recursos/J_history_simplificado.csv', 'w', newline='') as csvfile:
4       writer = csv.writer(csvfile)
5       writer.writerow(['Iteracion', 'J'])
6       slicedHistory = J_history[0:-1:100]
7       for i in range(1, len(slicedHistory) + 1):
8           writer.writerow([i*100, slicedHistory[i - 1]])
9   with open('./memoria/recursos/wb.csv', 'w', newline='') as csvfile:
10      writer = csv.writer(csvfile)
11      writer.writerow(X_features + ['b']),
12      writer.writerow(np.concatenate((w, [b])))
13
14   with open('./memoria/recursos/mean.csv', 'w', newline='') as csvfile:
15      writer = csv.writer(csvfile)
16      writer.writerow(X_features),
17      writer.writerow(mu)
18   with open('./memoria/recursos/sigma.csv', 'w', newline='') as csvfile:
19      writer = csv.writer(csvfile)
20      writer.writerow(X_features),
21      writer.writerow(sigma)
```

Figura 3.1: Código de la función 'write_results'