

Predicción del Precio de Viviendas en California.

Capítulo 1: Arquitectura de las Redes Neuronales

Código inicial para cargar los datos:

```
1 import numpy as np
2 from sklearn.datasets import fetch_california_housing
3 from sklearn.model_selection import train_test_split
4 from sklearn.preprocessing import StandardScaler
5
6 # Cargar el dataset
7 california = fetch_california_housing()
8 X = california.data
9 y = california.target
10 feature_names = california.feature_names
11
12 print('Descripción del dataset:')
13 print(california.DESCR[:1500])
14
15
```

Descripción del dataset:

.. _california_housing_dataset:

California Housing dataset

****Data Set Characteristics:****

:Number of Instances: 20640

:Number of Attributes: 8 numeric, predictive attributes and the target

:Attribute Information:

- MedInc median income in block group
- HouseAge median house age in block group
- AveRooms average number of rooms per household
- AveBedrms average number of bedrooms per household
- Population block group population
- AveOccup average number of household members
- Latitude block group latitude
- Longitude block group longitude

:Missing Attribute Values: None

This dataset was obtained from the StatLib repository.

https://www.dcc.fc.up.pt/~ltorgo/Regression/cal_housing.html

The target variable is the median house value for California districts, expressed in hundreds of thousands of dollars (\$100,000).

This dataset was derived from the 1990 U.S. census, using one row per census block group. A block group is the smallest geographical unit for which the U.S. Census Bureau publishes sample data (a block group typically has a population of 600 to 3,000 people).

A household is a group of people residing within a home. Since the average number of rooms and bedrooms in this dataset are provided per household, these columns may take surprisingly large values for block groups with few households and many empty houses, such as vacation resorts.

It can be downloaded/loaded

```
1 print(california.feature_names)
2 print(california.DESCR)
```

['MedInc', 'HouseAge', 'AveRooms', 'AveBedrms', 'Population', 'AveOccup', 'Latitude', 'Longitude']
.. _california_housing_dataset:

California Housing dataset

****Data Set Characteristics:****

:Number of Instances: 20640

:Number of Attributes: 8 numeric, predictive attributes and the target

```
:Attribute Information:
- MedInc          median income in block group
- HouseAge        median house age in block group
- AveRooms        average number of rooms per household
- AveBedrms       average number of bedrooms per household
- Population      block group population
- AveOccup        average number of household members
- Latitude        block group latitude
- Longitude       block group longitude
```

:Missing Attribute Values: None

This dataset was obtained from the StatLib repository.
https://www.dcc.fc.up.pt/~ltorgo/Regression/cal_housing.html

The target variable is the median house value for California districts, expressed in hundreds of thousands of dollars (\$100,000).

This dataset was derived from the 1990 U.S. census, using one row per census block group. A block group is the smallest geographical unit for which the U.S. Census Bureau publishes sample data (a block group typically has a population of 600 to 3,000 people).

A household is a group of people residing within a home. Since the average number of rooms and bedrooms in this dataset are provided per household, these columns may take surprisingly large values for block groups with few households and many empty houses, such as vacation resorts.

It can be downloaded/loaded using the
:func:`sklearn.datasets.fetch_california_housing` function.

.. rubric:: References

- Pace, R. Kelley and Ronald Barry, Sparse Spatial Autoregressions, Statistics and Probability Letters, 33 (1997) 291-297

1.

a) En este caso el numero de registros son **20640**

b) Cada observación tiene **8** características

c) El vector de características es:

$$\mathbf{X} = \begin{bmatrix} \text{MedInc} \\ \text{HouseAge} \\ \text{AveRooms} \\ \text{AveBedrms} \\ \text{Population} \\ \text{AveOccup} \\ \text{Latitude} \\ \text{Longitude} \end{bmatrix}$$

****2.****

MedInc: Ingreso medio de los hogares en el bloque censal (medido en decenas de miles de dolares).

HouseAge: Edad mediana de las viviendas en el bloque censal.

AveRooms: Numero promedio de habitaciones por hogar en el bloque censal.

AveBedrms: Numero promedio de dormitorios por hogar en el bloque censal.

Population: Numero total de personas que viven en el bloque censal.

AveOccup: Numero promedio de personas por hogar en el bloque censal.

Latitude: Latitud geografica del bloque censal, que indica su ubicación norte-sur en California.

Longitude: Longitud geografica del bloque censal, que indica su ubicación este-oeste en California.

3.

La variable objetivo es el valor medio de las viviendas en los distritos de California, expresado en cientos de miles de dolares (100.000 dólares).

****4.****

Este es un problema de regresión, dado que queremos obtener el valor de una vivienda en función de sus características ya mencionadas

Normalización de datos:

```
1 # Dividir en conjunto de entrenamiento y prueba
2 X_train, X_test, y_train, y_test = train_test_split(
3 X, y, test_size=0.2, random_state=42)
4
5 # Normalizar las características
6 scaler = StandardScaler()
7 X_train_scaled = scaler.fit_transform(X_train)
8 X_test_scaled = scaler.transform(X_test)
9 n_caracteristicas = X_train_scaled.shape[1]
```

d) La normalización de las características es importante porque garantiza que todas las variables de entrada tengan escalas comparables, evitando que aquellas con valores numéricos más grandes dominen el proceso de aprendizaje.

En el descenso del gradiente, esto hace que los pasos de actualización sean desbalanceados, provocando que el aprendizaje sea más lento o inestable. Al normalizar las características, el descenso del gradiente avanza de forma más uniforme, lo que permite una convergencia más rápida y estable del entrenamiento.

e) StandardScaler() se encarga de normalizar la distribución de los datos, tal que queden centrados en 0, y con dispersión 1.

$$x'_j = \frac{x_j - \mu_j}{\sigma_j}$$

Este proceso se hace para cada dato del conjunto de entrenamiento, de esta forma, se obtiene un conjunto con las mismas dimensiones.

f) La razón por la cual solo se aplica fit() los datos de entrenamiento y no a los de prueba, es justamente porque no queremos que en el entrenamiento se incluyan valores con los que vamos a realizar la prueba, de esta manera logramos un aislamiento entre estos 2 conjuntos. El transform() si se aplica sin problema.

```
1 def inicializar_parametros(n_caracteristicas):
2     '''
3     Inicializa los pesos y el sesgo del perceptrón.
4
5     Parámetros:
6     -----
7     n_caracteristicas : int
8     Número de características de entrada
9
10    Retorna:
11    -----
12    w : numpy array de forma (n_caracteristicas, 1)
13    b : float (escalar)
14    '''
15
16    # COMPLETE EL CÓDIGO AQUÍ
17    # Inicialice w con valores aleatorios pequeños (usar np.random.randn)
18    # Inicialice b en cero
19    w = 0.045*np.random.randn(n_caracteristicas, 1)
20    b = float(0)
21    return w, b
```

g) si los valores son cercanos a 0, se garantiza que no se alcanza una saturación muy rápida de la función de activación, puesto que permite que el gradiente aumente y lograr una correcta evolución. Además, iniciar todos los pesos con el mismo valor introduce un comportamiento de simetría indeseable a la hora de entrenar el modelo

h) w es un vector columna de en este caso 8 componentes, tal que sus dimensiones son (8,1)

i) en este caso que b sea 0 no afecta la evolución del gradiente, en contraste a se inicia con un valor grande que puede saturar la función de activación

```
1 #codigo para verificar las dimensiones del vector w
2 w= inicializar_parametros(8)[0]
3 print(w.shape)
```

(8, 1)

```
1 def propagacion_adelante(X, w, b):
2
```

```

3     '''
4     Calcula la suma ponderada para todas las observaciones.
5
6     Parámetros:
7     -----
8     X : numpy array de forma (m, n)
9     m observaciones, n características
10    w : numpy array de forma (n, 1)
11    Vector de pesos
12    b : float
13    Sesgo
14
15    Retorna:
16    -----
17    y_pred : numpy array de forma (m, 1)
18    Predicciones del modelo
19    '''
20
21    # COMPLETE EL CÓDIGO AQUÍ
22    # Calcule z = X @ w + b (producto matricial)
23
24    z = X @ w + b
25    y_pred = z
26    return y_pred

```

j) X tiene dimensiones (m,n) multiplicado por w que tiene dimensiones de (n,1) se traduce en: (m,n)*(n,1), la única forma que el producto estaría bien definido es para $X @ w$. esto resulta en un vector de dimensiones (m,1).

k) Como es un problema de regresión y no de clasificación, la función de activación recomendada es la identidad, que cumple $y_{pred} = z$. Esto porque solo queremos saber la respuesta de una variable a la variación de otros parámetros.

l)

```

1 w,b = inicializar_parametros(8)
2 y_pred = propagacion_adelante(X_train_scaled[0:10], w, b)
3 print(y_pred[0:10])
4 print(y_pred[5][0])
5 print("\n")
6
7 print(y_train[0:10])
8 print(y_train[5])

```

Respuesta predicha:

```

[[ 0.01885408]
 [ 0.02933514]
 [ 0.00275105]
 [ 0.03401993]
 [-0.0093871 ]
 [-0.02558981]
 [ 0.02014833]
 [ 0.00550068]
 [ 0.02225742]
 [-0.09104416]]
-0.025589805899004735

```

Respuesta real:

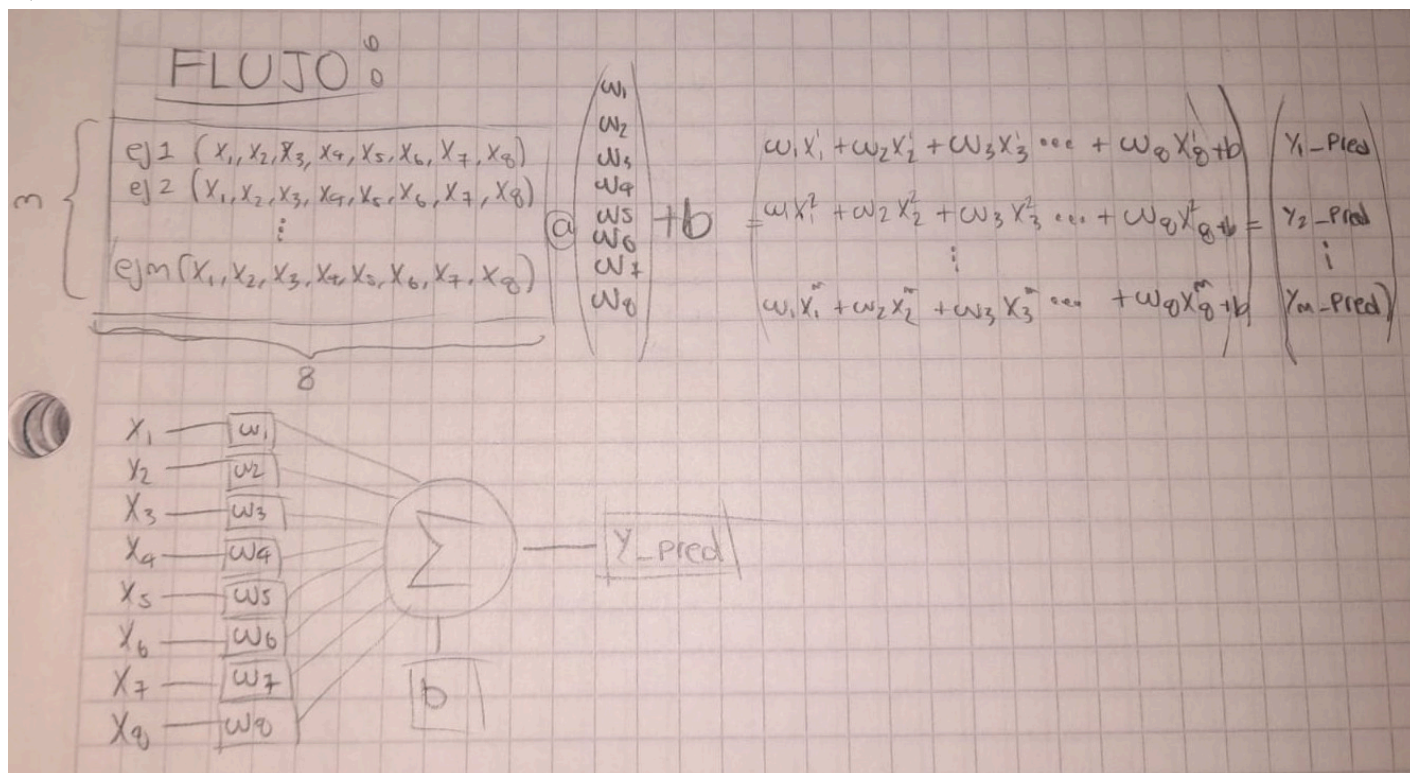
```

[1.03  3.821  1.726  0.934  0.965  2.648  1.573  5.00001 1.398
 3.156 ]
2.648

```

Estas son unas predicciones muy malas. Esto porque se hace el cálculo con los valores iniciales de los pesos y el sesgo, valores que se seleccionaron aleatoriamente con números muy pequeños además.

m)



```

1 def calcular_perdida(y_pred, y_real):
2     '''
3     Calcula el error cuadrático medio.
4
5     MSE = (1/m) * sum((y_pred - y_real)^2)
6     '''
7     # COMPLETE EL CÓDIGO AQUÍ
8
9     m = y_real.shape[0]
10    mse = (1/m)*np.sum((y_pred - y_real.reshape(-1,1))**2)
11    return mse
12
13

```

n) se prefiere usar la diferencia al cuadrado porque de esta manera se penaliza en gran medida, valores muy alejados y se preocupa poco de valores un poco mas cercanos.

o) la perdida inicial para 10 ejemplos es 6.8318

p) el valor absoluto no es diferenciable en 0, por otro lado, la diferencia al cuadrado es diferenciable en todo el dominio. Por eso el descenso del gradiente puede fallar en algun momento si se usa el MAE

q) que la funcion de perdida se estabilice puede significar varias cosas. Entre ellas que ya se llevo al minimo global, y que por mas que se avance en el entrenamiento no va a disminuir mas. Puede significar que se llegó a un minimo local, y no está bajando pero puede llegar a bajar mas.

$$\frac{\partial L}{\partial \omega} = \frac{\partial L}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial \omega}$$

$$1) \frac{\partial L}{\partial \hat{y}_i} = \frac{2}{m} \sum (\hat{y}_i - y_i) \quad , \quad \frac{\partial \hat{y}_i}{\partial \omega} = x_i^T \quad \frac{\partial L}{\partial \omega} = \left[\frac{2}{m} \sum (\hat{y}_i - y_i) \right] \cdot X_i^T$$

$$\boxed{\frac{\partial L}{\partial \omega} = \frac{2X^T}{m} \sum (\hat{y}_i - y_i)}$$

$$5) \frac{\partial L}{\partial b} = \frac{2}{m} \sum (\hat{y}_i - y_i) \quad \frac{\partial L}{\partial b} = \frac{\partial L}{\partial y_i} \frac{\partial y_i}{\partial b} = \frac{2}{m} \sum (\hat{y}_i - y_i) \cdot 1$$

$$\boxed{\frac{\partial L}{\partial b} = \frac{2}{m} \sum (\hat{y}_i - y_i)}$$

t) esto es por su definicion misma. Es un operador que evalua la funcion en ese punto y realizando la derivada en cada componente, arroja la direccion en la cual se tiene la distancia mas corta hasta el punto maximo. Que se traduce en la direccion de maximo ascenso, que en terminos de la analogia de la montaña, puede decirse que es la ruta que minimiza la distancia a la cumbre

u) el signo negativo significa que queremos bajar de la montaña, no subirla. Si el gradiente nos da la direccion de maximo ascenso su negativo (girar 180° ese vector) nos da el maximo descenso

```

1 def calcular_gradientes(X, y_pred, y_real):
2     '''
3     Calcula los gradientes de la pérdida respecto a w y b.
4
5     Parámetros:
6     -----
7     X : numpy array de forma (m, n)
8     y_pred : numpy array de forma (m, 1)
9     y_real : numpy array de forma (m, 1)
10
11     Retorna:
12     -----
13     dw : numpy array de forma (n, 1) - gradiente respecto a w
14     db : float - gradiente respecto a b
15     '''
16     m = X.shape[0]
17     error = y_pred - y_real.reshape(-1, 1)
18
19     # COMPLETE EL CÓDIGO AQUÍ
20     dw = (2/m)*X.T@error
21
22     db = (2/m)*np.sum(error)
23
24     return dw, db

```

1

w) dw y w deben tener la misma forma para que w tiene forma de (8,1) y dw al ser el producto de (8,m)@(m,1) el resultado es tambien (8,1). Es importante que tengan la misma forma para que la diferencia del siguiente paso se pueda ejecutar sin problema

x) si los valores calculados son iguales a los estimados, la diferencia será 0. Y entonces tanto dw como db será 0. Y esto significa que estamos en el minimo global, y no hay forma de descender mas

```

1 def actualizar_parametros(w, b, dw, db, learning_rate):
2     '''
3     Actualiza los pesos y sesgo usando gradiente descendente.
4
5     w_nuevo = w - learning_rate * dw
6     b_nuevo = b - learning_rate * db

```

```

7     '''
8     # COMPLETE EL CÓDIGO AQUÍ
9     w = w - learning_rate*dw
10    b = b - learning_rate*db
11
12    return w, b

```

z) este valor nos dice que tan grande o que tan pequeños serán los pasos que demos para "bajar la montaña". si es muy grande vamos a estar sobrepasando el punto optimo, no devolvemos y seguramente volvamos a rebasarlo. Si por otro lado es muy pequeño, tardará mucho tiempo en llegar al valor optimo

aa)

```

1 w, b = inicializar_parametros(8)
2 y_pred = propagacion_adelante(X_train_scaled[0:10], w, b)
3 mse = calcular_perdida(y_pred[0:10], y_train[0:10])
4 print(mse)
5 dw, db = calcular_gradientes(X_train_scaled[0:10], y_pred[0:10], y_train[0:10])
6 w,b = actualizar_parametros(w, b, dw, db, 0.01)
7
8 y_pred = propagacion_adelante(X_train_scaled[0:10], w, b)
9 mse = calcular_perdida(y_pred[0:10], y_train[0:10])
10 print(mse)

```

```

6.6689609053542656
6.175162824575837

```

el mse inicial fue 6.6689 luego de una iteracion es 6.1752

```

1 def entrenar_perceptron(X_train, y_train, learning_rate, epochs):
2     '''
3     Entrena un perceptrón para regresión.
4
5     Parámetros:
6     -----
7     X_train : datos de entrenamiento
8     y_train : etiquetas de entrenamiento
9     learning_rate : tasa de aprendizaje
10    epochs : número de épocas (iteraciones sobre todo el dataset)
11
12    Retorna:
13    -----
14    w, b : parámetros entrenados
15    historial_perdida : lista con la pérdida en cada época
16    '''
17    n_caracteristicas = X_train.shape[1]
18    w, b = inicializar_parametros(n_caracteristicas)
19    historial_perdida = []
20
21    for epoca in range(epochs):
22        # COMPLETE EL CÓDIGO AQUÍ
23        # 1. Propagación adelante
24        y_pred = propagacion_adelante(X_train_scaled, w, b)
25        # 2. Calcular pérdida y guardar en historial
26        perdida_actual = calcular_perdida(y_pred, y_train)
27        historial_perdida = historial_perdida + [perdida_actual]
28        # 3. Calcular gradientes
29        dw, db = calcular_gradientes(X_train_scaled, y_pred, y_train)
30        # 4. Actualizar parámetros
31        w,b = actualizar_parametros(w, b, dw, db, 0.01)
32
33        if epoca % 100 == 0:
34            print(f'Época {epoca}, Pérdida: {perdida_actual:.4f}')
35
36    return w, b, historial_perdida
37
38 w_entrenado, b_entrenado, historial_perdida = entrenar_perceptron(X_train, y_train, 0.01, 1000+1)
39 print(w_entrenado, "\n")
40 print(b_entrenado)

```

```

Época 0, Pérdida: 5.6872
Época 100, Pérdida: 0.7133
Época 200, Pérdida: 0.5966
Época 300, Pérdida: 0.5742
Época 400, Pérdida: 0.5592

```

```

Época 500, Pérdida: 0.5483
Época 600, Pérdida: 0.5403
Época 700, Pérdida: 0.5345
Época 800, Pérdida: 0.5303
Época 900, Pérdida: 0.5271
Época 1000, Pérdida: 0.5248
[[ 0.86086722]
 [ 0.1507098 ]
 [-0.25761414]
 [ 0.28514175]
 [ 0.00728149]
 [-0.0435653 ]
 [-0.67677527]
 [-0.64803947]]

```

2.0719469339614123

cc)

```

1 import matplotlib.pyplot as plt
2 plt.figure(figsize=(10, 6))
3 plt.plot(historial_perdida, color='blue', linewidth=2)
4 plt.title('Historial de Pérdida Vs Epoca', fontsize=16)
5 plt.xlabel('Epoca', fontsize=14)
6 plt.ylabel('Pérdida (MSE)', fontsize=14)
7 plt.grid(True, linestyle='--', alpha=0.7)
8 plt.xticks(fontsize=12)
9 plt.yticks(fontsize=12)
10 plt.show()

```



dd) es evidente que el modelo convergió puesto que se ve como desde la epoca 3000 no hay un descenso importante de la funcion de perdida. Se puede decir que con 4000 epocas de entrenamiento hubiese sido suficiente para tener un modelo util

ee)

```

1 plt.figure(figsize=(10, 6))
2
3 w_entrenado, b_entrenado, historial_perdida = entrenar_perceptron(X_train, y_train, 0.001, 1000+1)
4 plt.plot(historial_perdida, color='blue', linewidth=2, label='TA=0.001')
5
6 w_entrenado, b_entrenado, historial_perdida = entrenar_perceptron(X_train, y_train, 0.01, 1000+1)
7 plt.plot(historial_perdida, color='orange', linewidth=2, label='TA=0.01')
8
9 w_entrenado, b_entrenado, historial_perdida = entrenar_perceptron(X_train, y_train, 0.1, 1000+1)
10 plt.plot(historial_perdida, color='green', linewidth=2, label='TA=0.1')
11

```



```

12 w_entrenado, b_entrenado, historial_perdida = entrenar_perceptron(X_train, y_train, 1, 1000+1)
13 plt.plot(historial_perdida, color='red', linewidth=2, label='TA=1')
14
15 plt.title('Historial de Pérdida Vs Epoca', fontsize=16)
16 plt.xlabel('Epoca', fontsize=14)
17 plt.ylabel('Pérdida (MSE)', fontsize=14)
18 plt.grid(True, linestyle='--', alpha=0.7)
19 plt.xticks(fontsize=12)
20 plt.yticks(fontsize=12)
21 plt.legend()
22 plt.show()

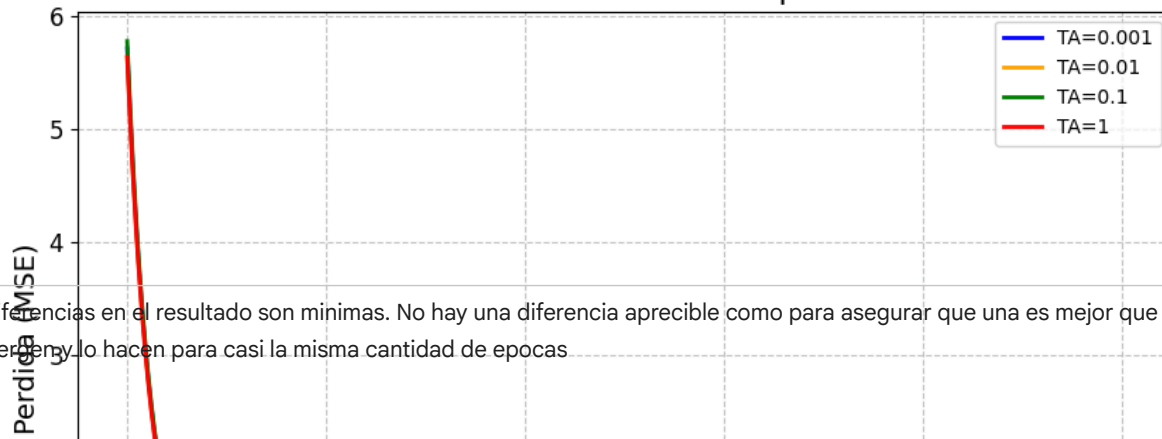
```

```

Época 0, Pérdida: 5.7131
Época 100, Pérdida: 0.7215
Época 200, Pérdida: 0.6016
Época 300, Pérdida: 0.5774
Época 400, Pérdida: 0.5612
Época 500, Pérdida: 0.5495
Época 600, Pérdida: 0.5410
Época 700, Pérdida: 0.5349
Época 800, Pérdida: 0.5304
Época 900, Pérdida: 0.5271
Época 1000, Pérdida: 0.5247
Época 0, Pérdida: 5.6403
Época 100, Pérdida: 0.7075
Época 200, Pérdida: 0.5928
Época 300, Pérdida: 0.5715
Época 400, Pérdida: 0.5573
Época 500, Pérdida: 0.5470
Época 600, Pérdida: 0.5395
Época 700, Pérdida: 0.5339
Época 800, Pérdida: 0.5299
Época 900, Pérdida: 0.5269
Época 1000, Pérdida: 0.5247
Época 0, Pérdida: 5.7773
Época 100, Pérdida: 0.7261
Época 200, Pérdida: 0.6056
Época 300, Pérdida: 0.5808
Época 400, Pérdida: 0.5640
Época 500, Pérdida: 0.5518
Época 600, Pérdida: 0.5429
Época 700, Pérdida: 0.5364
Época 800, Pérdida: 0.5317
Época 900, Pérdida: 0.5282
Época 1000, Pérdida: 0.5256
Época 0, Pérdida: 5.6353
Época 100, Pérdida: 0.6989
Época 200, Pérdida: 0.5864
Época 300, Pérdida: 0.5666
Época 400, Pérdida: 0.5535
Época 500, Pérdida: 0.5440
Época 600, Pérdida: 0.5371
Época 700, Pérdida: 0.5320
Época 800, Pérdida: 0.5284
Época 900, Pérdida: 0.5257
Época 1000, Pérdida: 0.5237

```

Historial de Pérdida Vs Epoca



```

1 w_entrenado, b_entrenado, historial_perdida = entrenar_perceptron(X_train, y_train, 0.001, 1000+1)
2 y_pred_test = X_test_scaled @ w_entrenado + b_entrenado
3 print(y_pred_test[-10:])
4 print()
5 print(y_test[-10:])

```

```

Época 0, Pérdida: 5.5566
Época 100, Pérdida: 0.7085

```

	200	400	600	800	1000
Época 200, Pérdida: 0.5966					
Época 300, Pérdida: 0.5744					
Época 400, Pérdida: 0.5593					
Época 500, Pérdida: 0.5484					
Época 600, Pérdida: 0.5404					
Época 700, Pérdida: 0.5345					
Época 800, Pérdida: 0.5302					
Época 900, Pérdida: 0.5271					
Época 1000, Pérdida: 0.5248					
[[4.3504323]					
[1.10118357]					
[1.95984691]					
[4.15532283]					
[5.70880786]					
[2.0211479]					
[2.07268185]					
[4.42535836]					
[1.30744207]					
[1.95643063]]					
[4.338 0.729 1.085 4.526 4.676 2.633 2.668 5.00001 0.723					
1.515]					

gg)

```
1 mse = np.mean((y_pred_test - y_test.reshape(-1,1))**2)
2 print(mse)
```

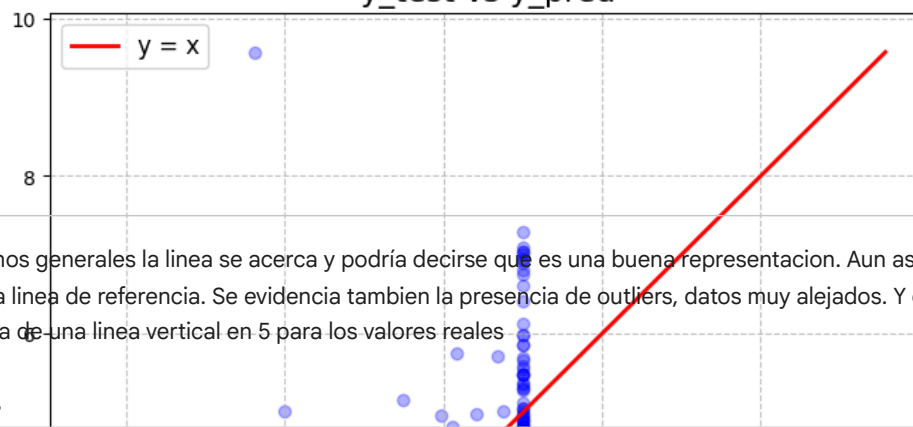
0.5548212650711372

el mse encontrado para el conjunto de prueba (0.5548) es casi igual al coseguido por el conjunto de entrenamiento (0.5247). Se puede concluir que el conjunto de entrenamiento generaliza apropiadamente el conjunto de prueba. Además, no hay muestras de sobre ajuste

hh

```
1 plt.figure(figsize=(8, 8))
2 plt.scatter(y_test, y_pred_test, alpha=0.3, color='blue')
3
4 # Línea ideal y = x
5 min_val = min(y_test.min(), y_pred_test.min())
6 max_val = max(y_test.max(), y_pred_test.max())
7 plt.plot([min_val, max_val], [min_val, max_val], color='red', linewidth=2, label='y = x')
8
9 plt.title('y_test vs y_pred', fontsize=16)
10 plt.xlabel('y_test (valores reales)', fontsize=14)
11 plt.ylabel('y_pred (valores predichos)', fontsize=14)
12 plt.grid(True, linestyle='--', alpha=0.7)
13 plt.legend(fontsize=12)
14 plt.show()
```

y_test vs y_pred



en terminos generales la linea se acerca y podria decirse que es una buena representacion. Aun asi hay una importante dispersion en torno a la linea de referencia. Se evidencia tambien la presencia de outliers, datos muy alejados. Y otro comportamiento raro es la existencia de una linea vertical en 5 para los valores reales

ii)

```
1 # Fórmula R²
2 SS_res = np.sum((y_test - y_pred_test.flatten())**2)
3 SS_tot = np.sum((y_test - np.mean(y_test))**2)
4 R2 = 1 - (SS_res / SS_tot)
5 print(R2)
```

0.5766044995998636

2

R2 representa una medida sobre que tanto una variable explica las variaciones de otra. En este caso dice que el 57% de los datos reales son explicados por los valores predichos. Lo cual es un valor medianmente eceptable. pero que se explica por lo mencionado anteriormente, algo de dispersion o outliers

jj) un perceptron simple unicamente puede modelar relaciones lineales entre las caracteristicas de entrada y un valor resultado (en este caso el precio de la vivienda)

kk) Si se tuviera un caso en que la variable dependiente no cumpliera una relacion lineal con las independientes, seria necesario incluir una capa oculta. Esto con el fin de incluir justamente un componente no lineal que explique este tipo de relaciones. Si se incluye una capa oculta con al menos una neurona mas vamos a tener una nueva ecuacion que al aplicarle determinada funcion de activacion permite doblar esa recta y acomodarse a esa no linealidad

ll) el teorema de aproximacion universal afirma que cualquier red neuronal con una unica capa oculta puede aproximar cualquier relacion. Pero este teorema no da informacion sobre cuantos perceptrones se deben usar, es decir, esto se cumple pero generalmetne requiere muchisimos perceptrones. Incluir varias capas ayuda a llegar al mismo resultado pero con una reduccion en el numero de neuronas requeridas

Fuentes