

## Práctica 5 (Train & Cross Validation)

### ¿Cómo he planteado la práctica?

Primero necesitamos los datos, en este caso leemos el .MAT 'ex5data1', y separamos los datos en X e Y, Xval e Yval, Xtest e Ytest.

Posteriormente hay que implementar la regresión lineal multivariable regularizada, que devolverá el coste y la gradiente.

Finalmente hay que dibujar la recta de resultados y la comparativa de los errores de aprendizaje.

## Código comentado y gráficas

### 1. Importar las librerías

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import scipy.optimize as opt
4
5 from scipy.io import loadmat
```

### 2. Obtener los datos del CSV:

```
1 data = loadmat('ex5data1.mat')
2
3 X = data['X']
4 Y = data['y'].ravel()
5 X_val = data['Xval']
6 Y_val = data['yval'].ravel()
7 X_test = data['Xtest']
8 Y_test = data['ytest'].ravel()
9
10 m = X.shape[0]
11 X = np.hstack([np.ones([m, 1]), X])
12 m_val = X_val.shape[0]
13 X_val = np.hstack([np.ones([m_val, 1]), X_val])
```

3. Ahora implementamos la regresión lineal multivariable, por un parte el coste, por otra parte la gradiente, y finalmente regularizamos ambas

```
1 def reg_lin_mul_var (theta, X, Y, lambd):
2
3     m = X.shape[0]
4     n = X.shape[1]
5     J_theta = 0
6     gradient = []
7
8     H = np.dot(X, theta)
9     diff = (H - Y)
10
11     # Coste
12     J_theta = (1 / (2 * m)) * np.sum(np.square(diff))
13
14     # Coste Regularizado
15     reg = lambd / (2 * m) * np.sum(np.square(theta[1:]))
16     J_theta += reg
17
18     # Gradiente
19     gradient = (1 / m) * np.sum(diff * X.T, axis=1)
20
21     # Gradiente Regularizada
22     reg = (lambd/m) * theta[1:]
23     gradient[1:] += reg
24
25
26     return (J_theta, gradient)
```

4. Y comprobamos los resultados

```
1 # PRUEBAS
2 theta = np.array([1, 1])
3 lambd = 1
4 cost, gradient = reg_lin_mul_var(theta, X, Y, lambd)
5 print("El coste es:\t\t{}\nLa gradiente es:\t{}".format(cost, gradient))
```

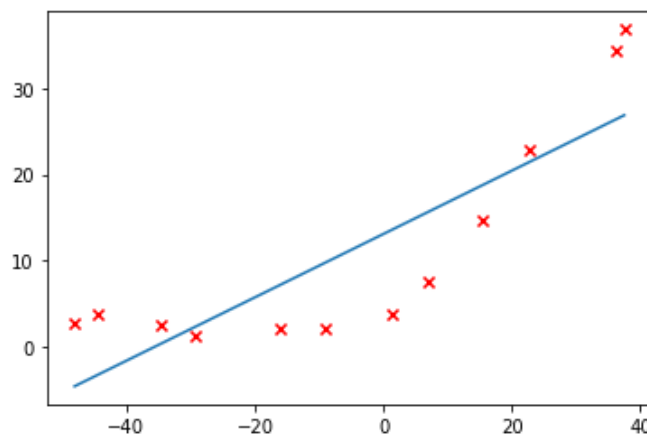
```
El coste es:          303.9931922202643
La gradiente es:      [-15.30301567 598.25074417]
```

5. Posteriormente minimizamos con la función MINIMIZE de Scipy, con  $\lambda = 0$ , el método escogido en este, y los siguientes puntos es 'TNC' → Truncated Newton

```
1 res = opt.minimize(  
2     fun = reg_lin_mul_var,  
3     x0 = theta,  
4     args=(X, Y, 0),  
5     method='TNC',  
6     jac=True  
7 )  
8  
9 theta_opt = res.x
```

6. Finalmente mostramos la recta de resultados, como en anteriores prácticas, nos quedamos con el mínimo y máximo, y predecimos la Y con las  $\theta_{OPTIMAS}$

```
1 plt.figure()  
2 plt.scatter(X[:, 1:], Y, marker='x', c='red')  
3  
4 X_min, X_max = np.min(X, axis=0), np.max(X, axis=0)  
5 Y_min = np.sum(theta_opt * X_min)  
6 Y_max = np.sum(theta_opt * X_max)  
7 X_min = X_min[1:]  
8 X_max = X_max[1:]  
9  
10 plt.plot([X_min, X_max], [Y_min, Y_max])  
11  
12 plt.show()
```



Podemos apreciar un alto sesgo (high bias), ya que los datos no se ajustan a la gráfica

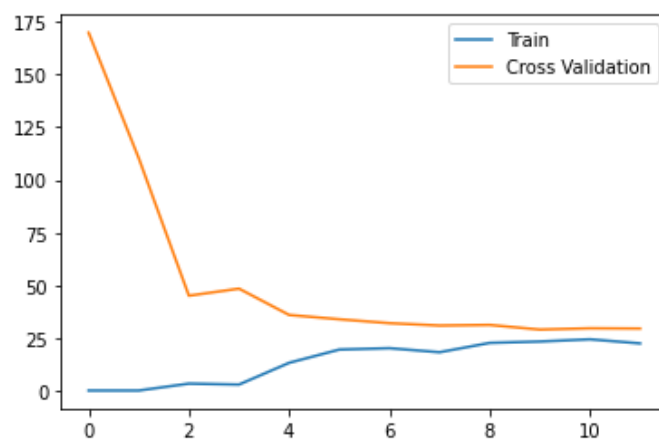
7. Ahora calculamos la tasa de error de aprendizaje, acumulando todos los ejemplos de  $X$  ( $X_i$ ), el error se calcula con la regresión lineal multivariable y con  $\lambda = 0$ . Esta función devuelve un array concatenado de todos los errores.

```
1 def learningErrors(theta, X, Y, X_val, Y_val, lambd = 0):
2     """
3     Returns the Train error vector and the Cross Validation error vector
4     """
5
6     trainError = []
7     cvError = []
8     m = X.shape[0]
9
10    for i in range(1, m + 1):
11
12        X_i = X[0:i]
13        Y_i = Y[0:i]
14
15        res = opt.minimize(
16            fun = reg_lin_mul_var,
17            x0 = theta,
18            args=(X_i, Y_i, lambd),
19            method='TNC',
20            jac=True
21        )
22
23        trainError.append(reg_lin_mul_var(res.x, X_i, Y_i, lambd)[0])
24        cvError.append(reg_lin_mul_var(res.x, X_val, Y_val, lambd)[0])
25
26    return np.array(trainError), np.array(cvError)
```

8. Finalmente, dibujamos el error de aprendizaje apoyándonos en la función anterior

```
1 def printLearningErrors(trainError, cvError):  
2     """  
3     Prints the Train error and the Cross Validation  
4     Doesn't return anything  
5     """  
6  
7     plt.figure()  
8  
9     r = range(0, len(trainError))  
10  
11    plt.plot(r, trainError)  
12    plt.plot(r, cvError)  
13    plt.legend(['Train', 'Cross Validation'])  
14  
15    plt.show()
```

```
1 trainError, cvError = learningErrors(theta, X, Y, X_val, Y_val)  
2 printLearningErrors(trainError, cvError)
```



Como habíamos comentado anteriormente, tenemos un alto sesgo, ya que el error de la Cross Validation es relativamente alto.

## Práctica 5 (Regresión polinomial)

### ¿Cómo he planteado la práctica?

Primero necesitamos los datos, en este caso leemos el .MAT 'ex5data1', y separamos los datos en X e Y, Xval e Yval, Xtest e Ytest.

Posteriormente hay que implementar una función que convierte los datos en vectores polinómicos  $x_1, x_2^2, \dots, x_p^p$ , y las funciones de normalización, una para normalizar un valor, y otra para normalizar a partir de unos valores  $\mu$  y  $\sigma$ .

En el siguiente paso, calculamos las predicciones y comprobamos si se ajusta a los datos.

Finalmente probamos varios parámetros  $\lambda$  y dibujamos el error que tiene respecto a su incremento, para así sacar un  $\lambda_{OPTIMO}$ .

## Código comentado y gráficas

### 1. Importar las librerías

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import scipy.optimize as opt
4
5 from scipy.io import loadmat
```

### 2. Obtener los datos del CSV:

```
1 data = loadmat('ex5data1.mat')
2
3 X = data['X']
4 Y = data['y'].ravel()
5 X_val = data['Xval']
6 Y_val = data['yval'].ravel()
7 X_test = data['Xtest']
8 Y_test = data['ytest'].ravel()
9
10 m = X.shape[0]
11 m_val = X_val.shape[0]
```

### 3. Creamos la función polinómica que nos devuelve un vector con $x_1, x_2^2, \dots, x_p^p$ :

```
1 def poly_data(X, p=1):
2     """
3     INPUT: Matix (m, 1) && p
4     OUTPUT: Matrix (m, p)
5     """
6     res = np.zeros((X.shape[0], p))
7
8     for i in range(0, p):
9         res[:, i] = X ** (i + 1)
10
11     return res
```



4. Implementamos una función de normalización a partir de un vector, con la siguiente fórmula:  $\frac{X-\mu}{\sigma}$   
Además, añade una columna de 1s.

```
1 def normalize(X):  
2     m = X.shape[0]  
3     mu = np.array([np.mean(X, axis = 0)])  
4     sigma = np.array([np.std(X, axis = 0)])  
5     X_norm = np.divide(np.subtract(X, mu), sigma)  
6  
7     X_norm = np.hstack([np.ones([m, 1]), X_norm])  
8  
9     return X_norm, mu, sigma
```

Por otro lado, implementamos una función de normalización a partir de unos valores  $\mu$  y  $\sigma$

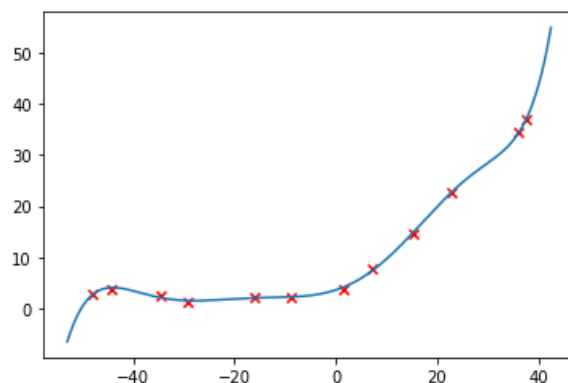
```
1 def normalize_with_values(X, mu, sigma):  
2     m = X.shape[0]  
3  
4     X_norm = np.divide(np.subtract(X, mu), sigma)  
5     X_norm = np.hstack([np.ones([m, 1]), X_norm])  
6  
7     return X_norm
```

5. Ahora usamos la función de MINIMIZE para obtener las  $\theta_{OPTIMAS}$ , para ello primero convertimos la X en polinómica y después la normalizamos (Recordar que también se añade la columna de 1s).

```
1 p = 8
2 theta = np.zeros(p + 1)
3 lambda = 0
4
5 X_poly = poly_data(X.ravel(), p)
6 X_poly_norm, mu, sigma = normalize(X_poly)
7
8 res = opt.minimize(
9     fun = reg_lin_mul_var,
10    x0 = theta,
11    args=(X_poly_norm, Y, lambda),
12    method='TNC',
13    jac=True
14 )
15
16 theta_opt = res.x
```

6. Y dibujamos la recta de resultados para comprobar si la polinomización se ajusta a los datos:

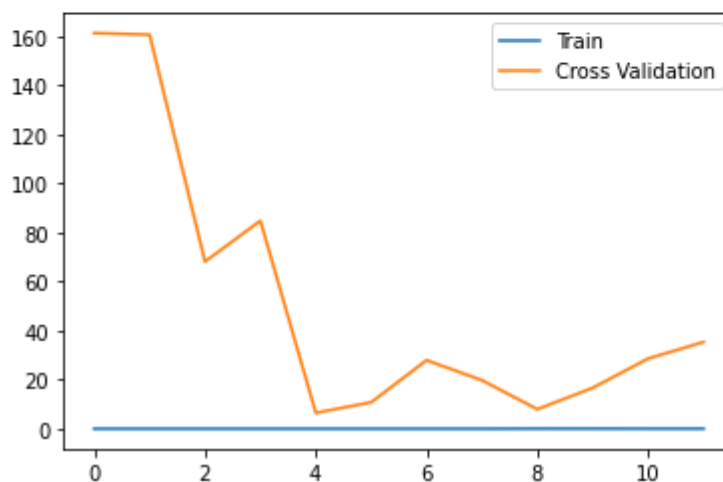
```
1 # Range of values -> -5 and +5 to show beyond the actual data
2 ranX = np.arange(
3     np.min(X) - 5,
4     np.max(X) + 5,
5     .05
6 )
7
8 # Need to normalize with mu and sigma
9 ranX_poly_norm = normalize_with_values(poly_data(ranX, p), mu, sigma)
10
11 # Predict values
12 Y_hat = np.dot(ranX_poly_norm, theta_opt)
13
14 # Plot the figure
15 plt.figure()
16 plt.scatter(X, Y, marker='x', c='red')
17 plt.plot(ranX, Y_hat)
18 plt.show()
```



Y podemos comprobar que sí, se ajusta a los datos, quizás demasiado, por ello vamos a dibujar la curva de aprendizaje.

7. Dibujamos la curva de aprendizaje, para ello necesitamos polinomializar y normalizar ambas  $X$ , la  $X_{TRAIN}$  y la  $X_{CV}$ .

```
1 X_poly_norm, mu, sigma = normalize(poly_data(X.ravel(), p))
2
3 # Need to normalize with mu and sigma
4 X_val_poly_norm = normalize_with_values(
5     poly_data(X_val.ravel(), p),
6     mu,
7     sigma
8 )
9
10 trainError, cvError = learningErrors(theta,
11                                     X_poly_norm,
12                                     Y,
13                                     X_val_poly_norm,
14                                     Y_val,
15                                     0
16 )
17
18 printLearningErrors(trainError, cvError)
```

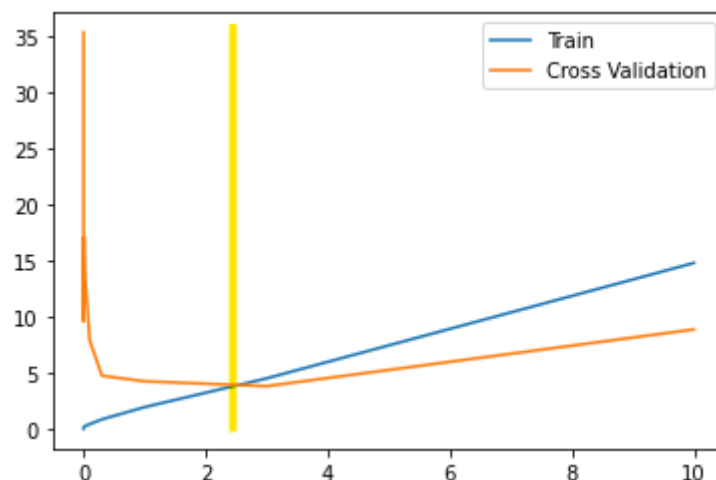


La curva de aprendizaje en el Train set es bajo y estable, lo cual está bien, por otro lado, la Cross Validation es baja pero irregular, por lo que necesitaríamos más datos para conseguir que bajase, o más iteraciones.

**NOTA: Recordar que en la curva de aprendizaje  $\lambda$  siempre debe ser 0**

8. Ahora dibujamos varias  $\lambda$  y su error, para conocer la  $\lambda_{OPTIMA}$ , al igual que antes, polinomizamos y normalizamos.

```
1 lambdas = [ 0, .001, .003, .01, .03, .1, .3, 1, 3, 10 ]
2 p = 8
3 trainError_vector = []
4 cvError_vector = []
5
6 plt.figure()
7
8 X_poly_norm, mu, sigma = normalize(poly_data(X.ravel(), p))
9
10 X_val_poly_norm = normalize_with_values(
11     poly_data(X_val.ravel(), p),
12     mu,
13     sigma
14 )
15
16
17 for lamdb in lambdas:
18     theta = np.zeros(p + 1)
19
20     res = opt.minimize(
21         fun = reg_lin_mul_var,
22         x0 = theta,
23         args=(X_poly_norm, Y, lamdb),
24         method='TNC',
25         jac=True,
26     )
27
28     theta_opt = res.x
29
30     trainError_vector.append(reg_lin_mul_var(theta_opt, X_poly_norm, Y, 0)[0])
31     cvError_vector.append(reg_lin_mul_var(theta_opt, X_val_poly_norm, Y_val, 0)[0])
32
33 plt.plot(lambdas, trainError_vector, label="Train")
34 plt.plot(lambdas, cvError_vector, label="Cross Validation")
35 plt.legend()
36 plt.show()
```



Podemos observar que el valor  $\lambda_{OPTIMO}$  está en el intervalo  $2 < \lambda < 3$

9. Finalmente, el valor más cercano a  $\lambda_{OPTIMA}$  que hemos probado es 3, por lo que comprobamos el error, y en este caso es de 3.5720, como pide el enunciado.

```
1  lambda_opt = 3
2  p = 8
3  theta = np.zeros(p + 1)
4
5  X_poly_norm, mu, sigma = normalize(poly_data(X.ravel()), p)
6
7  res = opt.minimize(
8      fun = reg_lin_mul_var,
9      x0 = theta,
10     args=(X_poly_norm, Y, lambda_opt),
11     method='TNC',
12     jac=True,
13 )
14
15 X_test_poly_norm = normalize_with_values(
16     poly_data(X_test.ravel(), p),
17     mu,
18     sigma
19 )
20
21 reg_lin_mul_var(res.x, X_test_poly_norm, Y_test, 0)[0]
```

3.572044856986457