

# Aprendizaje automático y minería de datos:

## Práctica 2

Jorge Rodríguez García y Gonzalo Sanz Lastra

### Código de la práctica:

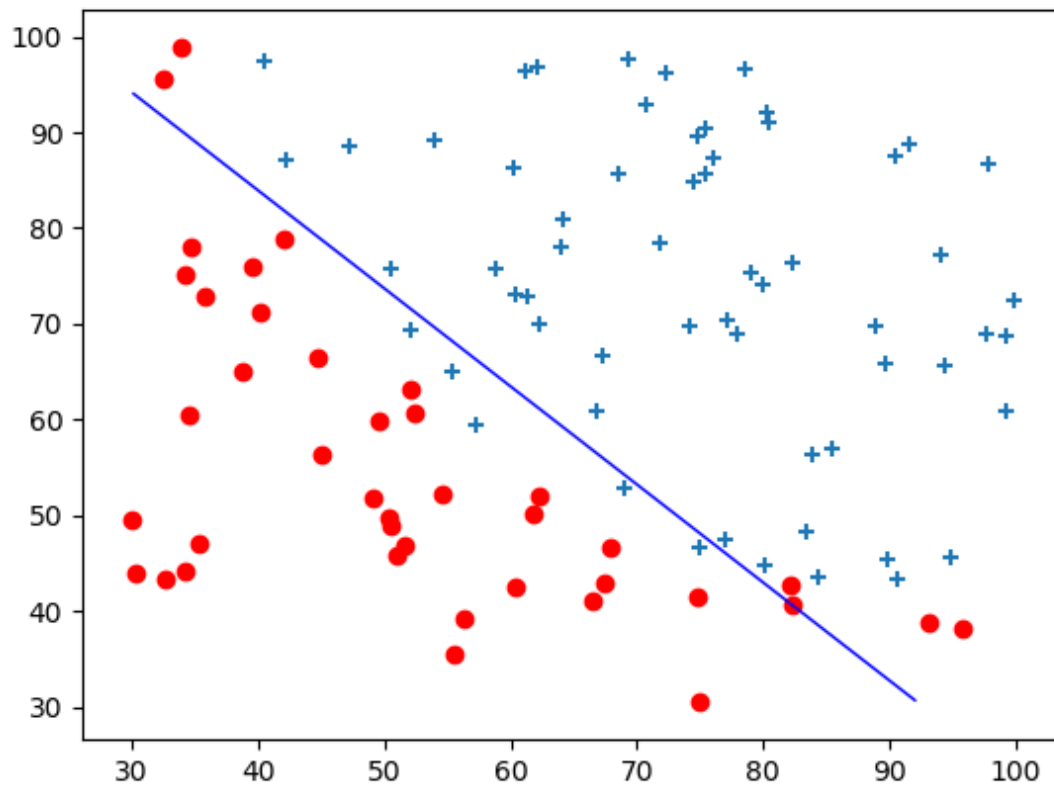
#### PARTE 1: REGRESIÓN LOGÍSTICA

```
1 import numpy as np
2 import scipy.optimize as opt
3 from pandas.io.parsers import read_csv # para leer .csv
4 from matplotlib import pyplot as plt # para dibujar las graficas
5
6 def load_csv(file_name):
7     """carga el fichero csv especificado y lo devuelve en un array de numpy"""
8     valores = read_csv(file_name, header=None).values
9
10    return valores.astype(float)
11
12 def graphics(X, Y, 0):
13     pinta_frontera_recta(X, Y, 0)
14
15     # Obtiene un vector con los índices de los ejemplos positivos
16     pos = np.where(Y == 1)
17     # Dibuja los ejemplos positivos
18     plt.scatter(X[pos, 1], X[pos, 2], marker= '+')
19
20     # Obtiene un vector con los índices de los ejemplos negativos
21     pos = np.where(Y == 0)
22     # Dibuja los ejemplos negativos
23     plt.scatter(X[pos, 1], X[pos, 2], c = 'red')
24
25     plt.savefig("frontera1.png")
26     plt.show()
27
28 def pinta_frontera_recta(X, Y, 0):
29     """pinta la recta que separa los datos entre los que cumplen el requisito y los que no"""
30     plt.figure()
31     x1_min, x1_max = X[:, 1].min(), X[:, 1].max()
32     x2_min, x2_max = X[:, 2].min(), X[:, 2].max()
33
34     xx1, xx2 = np.meshgrid(np.linspace(x1_min, x1_max),
35                             np.linspace(x2_min, x2_max)) # grid de cada columna de Xs
36
37     h = sigmoid(np.c_[np.ones((xx1.ravel().shape[0], 1)),
38                      xx1.ravel(), xx2.ravel()].dot(0)) # ravel las pone una tras otra
39
40     h = h.reshape(xx1.shape)
41
42     # el cuarto parámetro es el valor de z cuya frontera se
43     # quiere pintar, en este caso el punto medio entre 0 y 1 (los que cumplen y los que no)
44     plt.contour(xx1, xx2, h, [0.5], linewidths=1, colors='b')
45
46 # g(X*0t) = h(x)
47 def sigmoid(Z):
48     return 1/(1+np.exp(-Z))
49
```

```

50 def cost(O, X, Y):
51     """devuelve un valor de coste"""
52     return -((np.log(sigmoid(X.dot(O))))).T.dot(Y) + (np.log(1-sigmoid(X.dot(O))))).T.dot(1-Y))/X.shape[0]
53
54 def gradient(O, X, Y):
55     """la operacion que hace el gradiente por dentro -> devuelve un vector de valores"""
56     return (X.T.dot(sigmoid(X.dot(O))-Y))/X.shape[0]
57
58 def successPercentage(X, Y, O):
59     """determina el porcentaje de aciertos comparando los resultados estimados con los resultados reales"""
60     results = (sigmoid(X.dot(O.T)) >= 0.5)
61     results = (results == Y)
62     return results.sum()/results.shape[0]
63
64 def main():
65     valores = load_csv("ex2data1.csv")
66
67     X = valores[:, :-1] # matriz X, con todas las filas y todas las columnas menos la ultima (ys)
68     Y = valores[:, -1] # matriz Y, con todas las filas y la ultima columna
69
70     m = X.shape[0] # numero de muestras de entrenamiento
71     n = X.shape[1] + 1 # numero de variables x que influyen en el resultado y, mas la columna de 1s
72
73     X = np.hstack([np.ones([m, 1]), X])
74
75     O = np.zeros(n)
76
77     # por dentro hace la funcion de gradiente, usando nuestras funciones y variables
78     result = opt.fmin_tnc(func = cost, x0 = O, fprime = gradient, args=(X, Y))
79     O_opt = result[0] # entre todos los resultados devueltos, este ofrece las thetas optimas
80
81     success = successPercentage(X, Y, O_opt)
82     print("Porcentaje de acierto: " + str(success*100) + "%")
83
84     #GRAFICAS
85     graphics(X, Y, O_opt)
86
87 main()

```



## PARTE 2: REGULARIZACIÓN

```
1 import numpy as np
2 import scipy.optimize as opt
3 from pandas.io.parsers import read_csv # para leer .csv
4 from matplotlib import pyplot as plt # para dibujar las graficas
5 from sklearn import preprocessing # para polinomializar las Xs
6
7 def load_csv(file_name):
8     """carga el fichero csv especificado y lo devuelve en un array de numpy"""
9     valores = read_csv(file_name, header=None).values
10
11     return valores.astype(float)
12
13 def graphics(X, Y, O, poly):
14     plot_decisionboundary(X, Y, O, poly)
15
16     # Obtiene un vector con los índices de los ejemplos positivos
17     pos = np.where(Y == 1)
18     # Dibuja los ejemplos positivos
19     plt.scatter(X[pos, 0], X[pos, 1], marker= '+')
20
21     # Obtiene un vector con los índices de los ejemplos negativos
22     pos = np.where(Y == 0)
23     # Dibuja los ejemplos negativos
24     plt.scatter(X[pos, 0], X[pos, 1], c = 'red')
25
26     plt.savefig("frontera2.png")
27     plt.show()
28
29 def plot_decisionboundary(X, Y, theta, poly):
30     """pinta el polinomio que separa los datos entre los que cumplen el requisito y los que no"""
31     plt.figure()
32
33     x1_min, x1_max = X[:, 0].min(), X[:, 0].max()
34     x2_min, x2_max = X[:, 1].min(), X[:, 1].max()
35
36     xx1, xx2 = np.meshgrid(np.linspace(x1_min, x1_max),
37                             np.linspace(x2_min, x2_max)) # grid de cada columna de Xs
38
39     h = sigmoid(poly.fit_transform(np.c_[xx1.ravel(), xx2.ravel()]).dot(theta)) # ravel las pone una tras otra
40
41     h = h.reshape(xx1.shape)
42
43     plt.contour(xx1, xx2, h, [0.5], linewidths=1, colors='g')
44
45     # g(X*O) = h(x)
46     def sigmoid(Z):
47         return 1/(1+np.exp(-Z))
48
49     def cost(O, X, Y, l):
50
51         """devuelve un valor de coste"""
52         return (-((np.log(sigmoid(X.dot(O))))).T.dot(Y) + (np.log(1-sigmoid(X.dot(O))))).T.dot(1-Y))/X.shape[0]
53         + (1/(2*X.shape[0]))*(O[1:])**2).sum()) # igual que el coste anterior pero añadiendole esto ultimo para la regularizacion
54
55     def gradient(O, X, Y, l):
56         """la operacion que hace el gradiente por dentro -> devuelve un vector de valores"""
57         Aux0 = np.hstack([np.zeros([1]), O[1:]])) # sustituimos el primer valor de las thetas por 0 para que el termino independiente
58         # no se vea afectado por la regularizacion (lambda*0 = 0)
59         return (((X.T.dot(sigmoid(X.dot(O)))-Y)/X.shape[0])
60                 + (1/X.shape[0])*Aux0) # igual que el gradiente anterior pero añadiendole esto ultimo para la regularizacion
61
62     def successPercentage(X, Y, O):
63         """determina el porcentaje de aciertos comparando los resultados estimados con los resultados reales"""
64         results = (sigmoid(X.dot(O.T)) >= 0.5)
65         results = (results == Y)
66         return results.sum()/results.shape[0]
67
68 def main():
69     valores = load_csv("ex2data2.csv")
70
71     X = valores[:, :-1] # matriz X, con todas las filas y todas las columnas menos la ultima (ys)
72     Y = valores[:, -1] # matriz Y, con todas las filas y la ultima columna
73
74     m = X.shape[0] # numero de muestras de entrenamiento
75
76     polyGrade = 6 # a mayor grado, mayor ajuste
77     poly = preprocessing.PolynomialFeatures(polyGrade)
78     Xpoly = poly.fit_transform(X) # añade automaticamente la columna de 1s
79
80     n = Xpoly.shape[1] # numero de variables x que influyen en el resultado y, mas la columna de 1s
81
82     O = np.zeros(n)
83     l = 1 # cuanto mas se aproxime a 0, mas se ajustara el polinomio (menor regularizacion)
84
85     # por dentro hace la funcion de gradiente, usando nuestras funciones y variables
86     result = opt.fmin_tnc(func = cost, x0 = O, fprime = gradient, args=(Xpoly, Y, l))
87     O_opt = result[0] # entre todos los resultados devueltos, este ofrece las thetas optimas
88
89     success = successPercentage(Xpoly, Y, O_opt)
90     print("Porcentaje de acierto: " + str(success*100) + "%")
91
92     #GRAFICAS
93     graphics(X, Y, O_opt, poly)
94
95 main()
```

