

Práctica 2: Regresión logística

Rafael Herrera Troca
Rubén Ruperto Díaz

Aprendizaje Automático y Big Data
Doble Grado en Ingeniería Informática y Matemáticas

4 de septiembre de 2020

1. Introducción

En esta práctica hemos implementado una serie de funciones de coste y de gradiente para construir un modelo de regresión logística con el que predecir, en el primer apartado, si un alumno será admitido o no según sus notas en dos exámenes, y, en el segundo, si un microchip está defectuoso o no en función a los resultados de dos test de calidad. El resultado obtenido es una función que, dado un nuevo caso de prueba, devuelve un cierto valor que podemos considerar perteneciente a la clase positiva si es mayor o igual a 0.5 y a la negativa si no.

Hemos programado el código en el entorno *Spyder* con Python 3.7.

2. Procedimiento

Para la realización de la práctica, definimos una serie de funciones de coste y sus gradientes, que luego utilizaremos como parámetros de la función *fmin_tnc* de la librería *scipy.optimize*. Esta última función obtiene los valores de θ que minimizan la función de coste dada usando el gradiente que se le pasa.

Además, definimos la función sigmoide que, dado un cierto valor, lo transforma para que esté entre 0 y 1. Esta función es además la que utilizaremos para que el valor devuelto por el modelo logístico sea interpretable como se explicaba en la introducción: positivo si es mayor o igual que 0.5 y negativo si no lo es. La función sigmoide se define como sigue:

$$g(z) = \frac{1}{1 + e^{-z}}$$

2.1. Parte 1: Regresión logística

En este caso, definimos la función de coste según la fórmula que se muestra a continuación:

$$J(\theta) = \frac{1}{N}((\log(g(X\theta)))^T y + (\log(1 - g(X\theta)))^T (1 - y))$$

Esta función se presenta en su forma vectorizada, que es como la hemos implementado en el código, haciendo todas las operaciones de forma literal tal y como aparecen en la fórmula.

En cuanto a su gradiente, lo hemos implementado también en su forma vectorizada, y calculamos todos sus elementos a la vez, de forma que tras aplicar la función se obtiene un array con todos los elementos de dicho gradiente. La fórmula que hemos utilizado es la siguiente:

$$\frac{\partial(\theta)}{\partial\theta_j} = \frac{1}{N} X^T (g(X\theta) - y)$$

Una vez tenemos estas dos funciones, para el resto de la primera parte basta con usar *fmin_tnc* para obtener los valores de θ que hacen mínimo el coste. Estos valores serán los que definen nuestro modelo, de forma que, dado un nuevo dato X , la predicción para éste viene dada por $g(X\theta)$. Después, representamos los datos y la función obtenida en una gráfica y calculamos su porcentaje de aciertos. Los resultados obtenidos se explican en la sección correspondiente.

2.2. Parte 2: Regresión logística regularizada

En cuanto a la segunda parte, comenzamos nuevamente por definir la función de coste de forma vectorizada según la fórmula que sigue:

$$J(\theta) = -\frac{1}{N}((\log(g(X\theta)))^T y + (\log(1 - g(X\theta)))^T (1 - y)) + \frac{\lambda}{2N} \sum_{j=1}^n \theta_j^2$$

Para calcular su gradiente, también de forma vectorizada, seguimos la fórmula que se muestra a continuación:

$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{N} X^T (g(X\theta) - y) + \frac{\lambda}{N} \theta_j$$

Al igual que en la primera parte, calculamos todo el gradiente a la vez usando un array, no de forma iterativa. En este caso, hay que tener cuidado porque el término final de la fórmula sólo se suma cuando $j \neq 0$. Para hacer esto, sumamos un array que tiene en cada posición el valor de θ_j correspondiente y, en la primera, un 0.

Además, en esta parte, los datos están repartidos más o menos en dos círculos concéntricos, por lo que una recta no es suficiente para separarlos. Sin embargo, como sólo tenemos dos variables por cada dato, no hay suficientes grados de libertad para construir una función que lo estime con suficiente fiabilidad. Para solucionarlo, usamos la clase *PolynomialFeatures* de la librería *sklearn.preprocessing*, que dados unos datos y un grado, construye todas las posibles potencias de estos datos hasta dicho grado. Así, al usarla con grado 6, nuestras 2 variables se transforman en 28, que nos proporcionan la flexibilidad suficiente como para que la función de decisión se adapte correctamente a la distribución de los datos.

Una vez explicado esto, solo resta usar nuevamente *fmin_tnc* para calcular el valor de θ que minimiza la función de coste, representar el resultado y calcular el porcentaje de aciertos.

Finalmente, estudiamos más en profundidad cómo varía el modelo construido según el valor de λ . Para ello, hemos puesto este parámetro como uno más de los que recibe la función de coste y el gradiente. Así, usando un bucle hemos ido probando distintos valores con una escala logarítmica de base 10, tomándolos entre 10^{-10} y 10^2 . Para cada valor, hemos representado el modelo obtenido y almacenado su porcentaje de acierto. Finalmente, hemos representado todos los porcentajes de acierto en una gráfica. Se pueden encontrar detalles específicos sobre los resultados en la siguiente sección.

3. Resultados

3.1. Parte 1: Regresión logística

En el primer caso, buscamos un modelo capaz de predecir si un alumno será admitido a una universidad en función de sus resultados en dos exámenes. La función de coste alcanza un valor mínimo de $J(\theta) = 0,203$ para $\theta = (-25,161, 0,206, 0,201)$, por tanto, la recta $X\theta = 0$, que define la frontera de decisión entre los puntos que se predicen como admitidos o como rechazados, viene dada por la ecuación $0,206x_1 + 0,201x_2 - 25,161 = 0$.

Esta recta además coincide con la curva de nivel de altura 0,5 de la función sigmoide $g(X\theta)$, siendo $\theta = (-25,161, 0,206, 0,201)$ y $X = (1, x_1, x_2)$. Esto se debe a que $g(X\theta) = 0,5 \iff X\theta = 0$, donde $X\theta = 0$ es precisamente la ecuación de la frontera de decisión. Por tanto, podemos utilizar la función sigmoide como un clasificador, interpretando que un alumno será admitido si, al introducir sus calificaciones en la función, el resultado es $\geq 0,5$ y rechazado en caso contrario. Utilizando este procedimiento, nuestro modelo obtiene un 89 % de éxito al clasificar los datos de entrenamiento.

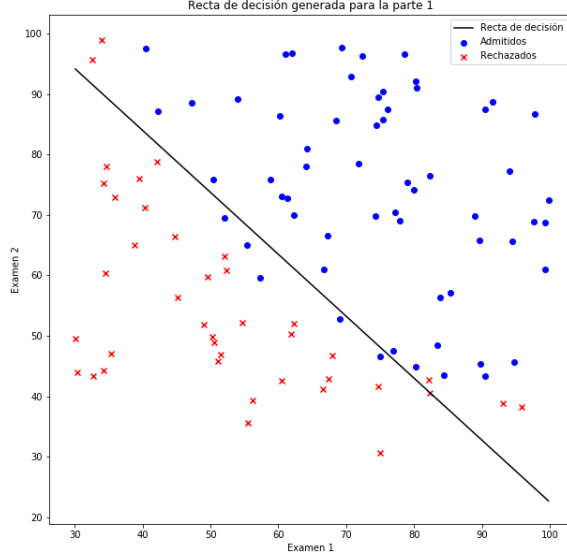


Figura 1: Recta de regresión para la primera parte

3.2. Parte 2: Regresión logística regularizada

En este segundo apartado debemos predecir si un microchip está defectuoso o no en función de los resultados de dos test. Sin embargo, como a la vista de los gráficos los datos no son separables mediante una recta, la ecuación que define la frontera de decisión tendrá que tener términos no lineales que combinen las variables x_1 y x_2 . De este modo, optimizamos un vector θ más grande, concretamente de 28 componentes, para conseguir una curva de hasta grado 6, y añadimos un término de regularización que depende de λ a la función de coste $J(\theta)$.

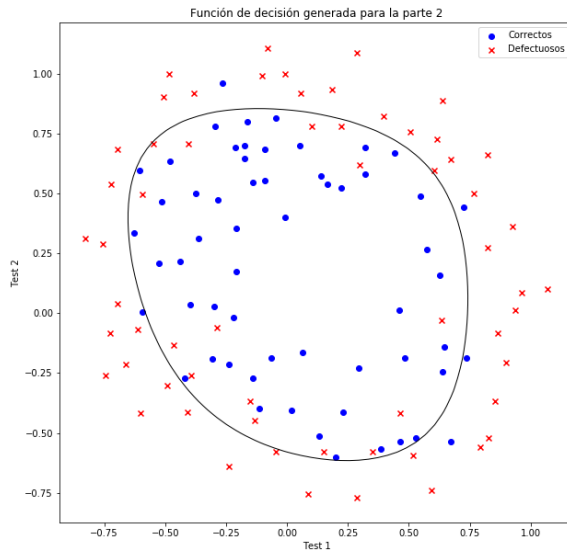


Figura 2: Frontera de decisión para $\lambda = 1$

Tras minimizar la función de coste regularizada con $\lambda = 1$, obtenemos una frontera de decisión que representamos gráficamente mediante la función *contour*, basándonos en que la curva que buscamos es la curva de nivel de altura 0,5 de la función sigmoide, fijado el parámetro θ que minimiza la función de coste regularizada.

Una vez estudiado el caso $\lambda = 1$, hemos probado el comportamiento de la regresión al dar a λ los distintos valores que hay en el conjunto $\{10^i \mid i \in [-10, 2], i \in \mathbb{Z}\}$. Hemos realizado una gráfica que muestra el porcentaje de casos de entrenamiento clasificados correctamente en función de los valores dados a λ , siendo importante notar que la escala del eje horizontal, donde se muestran los valores de λ , es logarítmica, por lo que los valores representados son sucesivas potencias de 10.

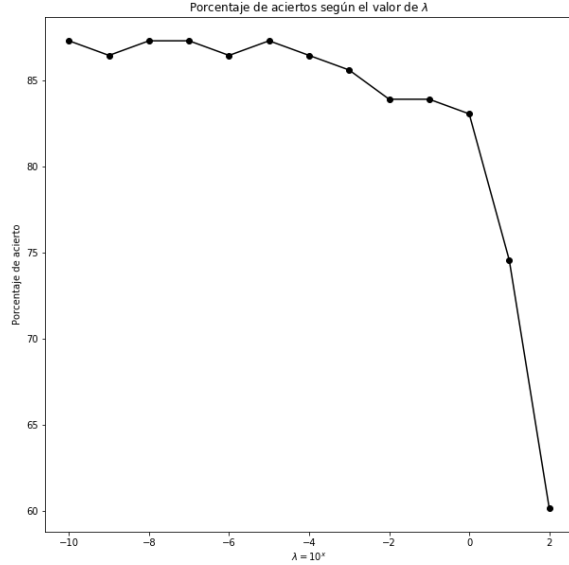


Figura 3: Comparativa de éxito al clasificar los datos de entrenamiento según λ

Vemos que si $\lambda \in \{0,01, 0,1, 1\}$ se tiene un porcentaje de aciertos entre 80 y 85 %. Para valores mayores ($\lambda \in \{10, 100\}$) la gráfica se desploma, indicando que la curva no se puede ajustar correctamente. Esto se debe a que aumentar λ tiene como consecuencia la disminución del valor de las componentes de θ que acompañan a x_1 , x_2 y sus combinaciones, y por tanto se restringe la forma de la curva y se reduce su tamaño.

Por otro lado, valores de $\lambda \leq 0,001$ obtienen tasas de acierto de más del 85 % con los datos de entrenamiento, pero esto en realidad no es bueno, ya que se debe a un sobreaprendizaje del modelo. La curva de decisión se ajusta demasiado a los datos de entrenamiento, permitiendo más aciertos entre esos datos, pero disminuyendo el éxito para predicciones con datos nuevos.

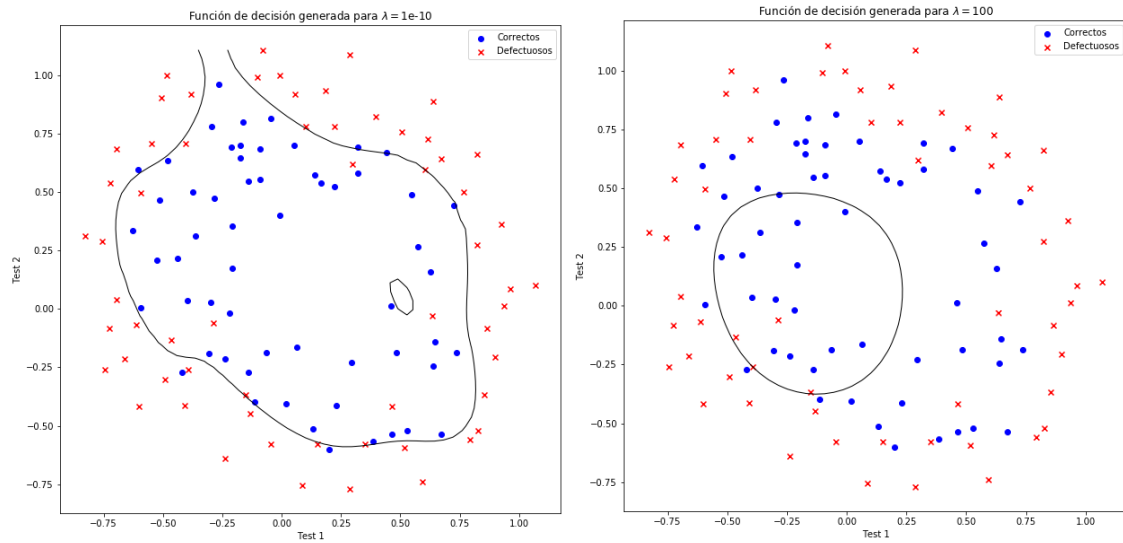


Figura 4: A la izquierda, la curva correspondiente a $\lambda = 10^{-10}$, con una forma excesivamente adaptada a los datos de entrenamiento. A la derecha, la curva para $\lambda = 100$ es demasiado pequeña como para englobar muchos de los puntos correspondientes a chips correctos.

4. Conclusiones

Esta práctica nos ha permitido comprender con gran nivel de detalle el funcionamiento de la regresión logística como método para construir modelos de predicción. Gracias a los gráficos realizados hemos podido visualizar el modelo obtenido y ver cómo se ajusta más o menos a los datos, sobre todo en la segunda parte al ir variando λ .

Además, hemos aprendido a utilizar la función *fmin_tnc* para minimizar el coste de $J(\theta)$ una vez se tiene la función de coste y su gradiente, lo cual es mucho más cómodo que implementar el proceso como se hacía en la práctica anterior, y a utilizar la clase *PolynomialFeatures* para aumentar el número variables que tiene cada dato.

5. Anexo: Código

A continuación se ofrece el código escrito para resolver los ejercicios propuestos en la práctica.

```
1  '''
2  Práctica 2
3
4  Rubén Ruperto Díaz y Rafael Herrera Troca
5  '''
6
7  import os
8  import numpy as np
9  import pandas as pd
10 import matplotlib.pyplot as plt
11 from scipy.optimize import fmin_tnc
12 from sklearn.preprocessing import PolynomialFeatures
13
14
15 os.chdir("./resources")
16
17
18 # Función sigmoide
19 def sigmoide(z):
20     return 1 / (1 + np.exp(-z))
21
22 def costeP1(theta, X, Y):
23     gXTheta = sigmoide(np.dot(X, theta))
24     factor = np.dot(np.log(gXTheta).T, Y) + np.dot(np.log(1 -
25 gXTheta).T, 1-Y)
26     return -1 / len(Y) * factor
27
28 def gradienteP1(theta, X, Y):
29     gXTheta = sigmoide(np.dot(X, theta))
30     return 1 / len(Y) * np.dot(X.T, gXTheta-Y)
31
32 def acierto(X, Y, theta):
33     gXTheta = sigmoide(np.dot(X, theta))
34     resultados = [((gXTheta >= 0.5) & (Y == 1)) | ((gXTheta < 0.5)
35 & (Y == 0))]
36     return np.count_nonzero(resultados) / len(Y)
37
38 def costeP2(theta, X, Y, lmb=1):
39     gXTheta = sigmoide(np.dot(X, theta))
40     factor = np.dot(np.log(gXTheta).T, Y) + np.dot(np.log(1 -
41 gXTheta).T, 1-Y)
42     return -1 / len(Y) * factor + lmb / (2 * len(Y)) * np.sum(theta
43 **2)
44
45 def gradienteP2(theta, X, Y, lmb=1):
46     gXTheta = sigmoide(np.dot(X, theta))
47     thetaJ = np.concatenate(([0], theta[1:]))
```

```

44     return 1 / len(Y) * np.dot(X.T, gXTheta-Y) + lmb / len(Y) *
        thetaJ
45
46
47 ## Parte 1: Regresión con una variable
48
49
50 # Leemos los datos
51 data = pd.read_csv("ex2data1.csv", names=['exam1', 'exam2', '
        resultado'])
52 X = np.array([data['exam1'], data['exam2']]).T
53 Y = np.array(data['resultado'])
54
55 # Representamos los datos en una gráfica
56 adm = np.where(Y == 1)
57 rech = np.where(Y == 0)
58 plt.figure(figsize=(10,10))
59 plt.scatter(X[adm, 0], X[adm, 1], marker='o', c='blue', label="
        Admitidos")
60 plt.scatter(X[rech, 0], X[rech, 1], marker='x', c='red', label="
        Rechazados")
61 plt.title("Datos proporcionados para la parte 1")
62 plt.xlabel("Examen 1")
63 plt.ylabel("Examen 2")
64 plt.legend(loc="upper right")
65 plt.savefig("p1Datos.png")
66 plt.show()
67
68 # Calculamos el valor óptimo de theta
69 X2 = np.hstack((np.array([np.ones(len(Y))]).T,X))
70 theta0 = np.zeros(np.shape(X2)[1])
71 result = fmin_tnc(func=costeP1, x0=theta0, fprime=gradienteP1, args
        =(X2, Y))[0]
72
73 # Representamos la recta en una gráfica junto a los datos
74 part = np.linspace(min(X[:,0]), max(X[:,0]), 1000)
75 plt.figure(figsize=(10,10))
76 plt.scatter(X[adm, 0], X[adm, 1], marker='o', c='blue', label="
        Admitidos")
77 plt.scatter(X[rech, 0], X[rech, 1], marker='x', c='red', label="
        Rechazados")
78 plt.plot(part, -(result[0]+result[1]*part)/result[2], 'k', label="
        Recta de decisión")
79 plt.title("Recta de decisión generada para la parte 1")
80 plt.xlabel("Examen 1")
81 plt.ylabel("Examen 2")
82 plt.legend(loc="upper right")
83 plt.savefig("p1Recta.png")
84 plt.show()
85
86 # Comprobamos el porcentaje de aciertos de la recta obtenida
87 porc_ac = acierto(X2, Y, result)
88 print("Se tiene un " + str(porc_ac*100) + "% de aciertos")
89
90

```



```

91 ## Parte 2
92
93
94 # Leemos los datos
95 data = pd.read_csv("ex2data2.csv", names=['test1', 'test2', '
    resultado'])
96 X = np.array([data['test1'], data['test2']]).T
97 Y = np.array(data['resultado'])
98
99 # Representamos los datos en una gráfica
100 corr = np.where(Y == 1)
101 defe = np.where(Y == 0)
102 plt.figure(figsize=(10,10))
103 plt.scatter(X[corr, 0], X[corr, 1], marker='o', c='blue', label="
    Correctos")
104 plt.scatter(X[defe, 0], X[defe, 1], marker='x', c='red', label="
    Defectuosos")
105 plt.title("Datos proporcionados para la parte 2")
106 plt.xlabel("Test 1")
107 plt.ylabel("Test 2")
108 plt.legend(loc="upper right")
109 plt.savefig("p2Datos.png")
110 plt.show()
111
112 # Aplicamos PolynomialFeatures para mapear los atributos y obtener
113 # un mejor ajuste a los ejemplos de entrenamiento
114 X2 = PolynomialFeatures(6).fit_transform(X)
115
116 # Calculamos el valor óptimo de theta
117 theta0 = np.zeros(np.shape(X2)[1])
118 result = fmin_tnc(func=costeP2, x0=theta0, fprime=gradienteP2, args
    =(X2, Y))[0]
119
120 # Representamos el resultado obtenido
121 xx1, xx2 = np.meshgrid(np.linspace(min(X[:,0]),max(X[:,0])), np.
    linspace(min(X[:,1]),max(X[:,1])))
122 mesh = PolynomialFeatures(6).fit_transform(np.c_[xx1.ravel(),xx2.
    ravel()])
123 h = sigmoide(np.dot(mesh,result))
124 h = h.reshape(xx1.shape)
125 plt.figure(figsize=(10,10))
126 plt.scatter(X[corr, 0], X[corr, 1], marker='o', c='blue', label="
    Correctos")
127 plt.scatter(X[defe, 0], X[defe, 1], marker='x', c='red', label="
    Defectuosos")
128 plt.contour(xx1, xx2, h, [0.5], linewidths=1, colors='k')
129 plt.title("Función de decisión generada para la parte 2")
130 plt.xlabel("Test 1")
131 plt.ylabel("Test 2")
132 plt.legend(loc="upper right")
133 plt.savefig("p2FuncionL1.png")
134 plt.show()
135
136 # Comprobamos el porcentaje de aciertos de la recta obtenida
137 porc_ac = acierto(X2, Y, result)

```

```

138 print("Se tiene un " + str(porc_ac*100) + "% de aciertos")
139
140 # Repetimos el proceso anterior con distintos valores de lambda
141 porc_ac = []
142 for lmb in range(-10, 3):
143     theta0 = np.zeros(np.shape(X2)[1])
144     result = fmin_tnc(func=costeP2, x0=theta0, fprime=gradienteP2,
145                       args=(X2, Y, 10**lmb))[0]
146
147     xx1, xx2 = np.meshgrid(np.linspace(min(X[:,0]),max(X[:,0])), np
148                             .linspace(min(X[:,1]),max(X[:,1])))
149     mesh = PolynomialFeatures(6).fit_transform(np.c_[xx1.ravel(),
150                                                     xx2.ravel()])
151     h = sigmoide(np.dot(mesh,result))
152     h = h.reshape(xx1.shape)
153     plt.figure(figsize=(10,10))
154     plt.scatter(X[corr, 0], X[corr, 1], marker='o', c='blue', label=
155                 "Correctos")
156     plt.scatter(X[defe, 0], X[defe, 1], marker='x', c='red', label=
157                 "Defectuosos")
158     plt.contour(xx1, xx2, h, [0.5], linewidths=1, colors='k')
159     plt.title(r"Función de decisión generada para $\lambda = $" +
160              str(10**lmb))
161     plt.xlabel("Test 1")
162     plt.ylabel("Test 2")
163     plt.legend(loc="upper right")
164     plt.savefig("p2FuncionVL" + str(lmb) + ".png")
165     plt.show()
166
167     porc_ac.append(acierto(X2, Y, result)*100)
168     print("Se tiene un " + str(porc_ac[-1]) + "% de aciertos para
169           lambda " + str(10**lmb))
170
171 # Representamos el porcentaje de aciertos según el valor de lambda
172 plt.figure(figsize=(10,10))
173 plt.plot(range(-10,3), porc_ac, 'ko-')
174 plt.title(r"Porcentaje de aciertos según el valor de $\lambda$")
175 plt.xlabel(r"$\lambda = 10^{\{x\}}$")
176 plt.ylabel("Porcentaje de acierto")
177 plt.savefig("AcXLambda.png")
178 plt.show()

```