

# Práctica 4:

## Entrenamiento de redes neuronales

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 las funciones de coste y gradiente de una red neuronal, así como la función de retro-propagación para su entrenamiento. Posteriormente, hemos estudiado como cambian los resultados de la red al variar ciertos parámetros del entrenamiento.

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

# 2. Procedimiento

Comenzamos por explicar las funciones más sencillas que hemos definido, mientras que las funciones más complejas para el cálculo del coste y su gradiente se desarrollarán en secciones posteriores.

En primer lugar, tenemos la función *sigmoide*, que sirve para transformar el resultado devuelto por una capa de la red en un valor en el intervalo  $[0, 1]$ , y que viene definida según la fórmula:

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

La segunda función definida es la derivada de la sigmoide, *diffSigmoide*, cuya definición para una capa  $l$  de la red, siendo  $a^{(l)}$  la entrada de dicha capa y  $z^{(l)}$  la salida de la capa  $l - 1$ , viene dada por la fórmula:

$$g'(z^{(l)}) = a^{(l)}(1 - a^{(l)})$$

Si tenemos en cuenta que  $a^{(l)} = [1 \sim g(z^{(l)})]$ , esto es, calcular la sigmoide del resultado devuelto por la capa anterior y poner una columna de unos delante (normalmente referida como  $a_0^{(l)}$ ), podemos expresar la derivada de la sigmoide según esta otra fórmula:

$$g'(z^{(l)}) = [1 \sim g(z^{(l)})](1 - [1 \sim g(z^{(l)})])$$

No obstante, por ser la implementación mucho más sencilla de este modo, en el código se ha utilizado la primera fórmula.

También hemos definido las funciones *applyLayer* y *applyNet* que, como su nombre indica, aplican una capa de la red dada la entrada y la matriz de pesos correspondientes, o la red entera dadas la entrada y un array con las matrices de pesos correspondientes a cada capa. Además, *applyNet* no sólo devuelve el resultado de aplicar la red, también devuelve la entrada utilizada en cada una de las capas, ya que necesitábamos este dato para calcular el gradiente de la función de coste.

Finalmente, también hemos definido las funciones *randomWeights*, que inicializa una matriz de pesos aleatorios para las dimensiones indicadas, y *acierto*, que calcula el porcentaje de acierto dadas la respuesta y las etiquetas correctas.

A continuación se explica con mayor grado de detalle el desarrollo de la función de coste y su gradiente.

## 2.1. Función de coste

Para la función de coste hemos aplicado la fórmula correspondiente manteniendo la generalidad todo lo posible. Dicha fórmula viene dada por:

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K \left[ -y_k^{(i)} \log((h_\theta(x^{(i)}))_k) - (1 - y_k^{(i)}) \log(1 - (h_\theta(x^{(i)}))_k) \right] + \frac{\lambda}{2N} \sum_{l=1}^L \left[ \sum_{i,j>0} (\Theta_{i,j}^{(l)})^2 \right]$$

En esta fórmula,  $N$  indica el número de casos,  $K$ , el número de etiquetas y  $L$ , el número de matrices de pesos que tiene la red. Asimismo,  $y$  es la salida esperada de la red neuronal,  $x$  es la entrada,  $h_\theta(x)$  es el resultado devuelto,  $\lambda$  es el parámetro de regularización y  $\Theta$  es la matriz de pesos de una cierta capa. Además, el sumatorio que aparece al final del término de regularización expresa que se debe sumar el cuadrado de todos los términos de cada matriz de pesos teniendo en cuenta no sumar la primera columna (recordemos que el índice de las filas empieza en uno siempre).

En cualquier caso, para aprovechar la potencia del cálculo matricial que ofrece Python, hemos vectorizado la función de coste, evitando así tener que hacer el cálculo de la primera parte para cada uno de los casos de prueba. La fórmula que hemos utilizado es la siguiente:

$$J(\theta) = \frac{1}{N} \sum_{i,j} [y \odot \log(h_\theta(x)) - (1 \ominus y) \odot \log(1 - h_\theta(x))]_{i,j} + \frac{\lambda}{2N} \sum_{l=1}^L \left[ \sum_{i,j>0} (\Theta_{i,j}^{(l)})^2 \right]$$

En esta fórmula, las letras vuelven a significar lo mismo que en la fórmula anterior, el término  $(1 \ominus y)$  representa la matriz  $M$  dada por  $M_{i,j} = 1 - y_{i,j}$ , y el símbolo  $\odot$  indica el producto elemento a elemento de dos matrices o producto de Hadamard.

Tras la construcción de la función, comprobamos con las matrices de pesos dadas como ejemplo que el resultado era el correcto.

## 2.2. Gradiente de la función de coste

Para el cálculo del gradiente implementamos la función *gradiente*, que lo calcula, nuevamente, de la forma más genérica posible para que se pueda utilizar para una red neuronal cualquiera sin tener que adaptarlo. Recibe como argumentos la salida

esperada para la red neuronal, la entrada de cada capa, el resultado de aplicar la red, la lista de matrices de pesos y el parámetro de regularización, y devuelve la lista de matrices de gradiente.

Para obtener este resultado vamos calculando la matriz del gradiente correspondiente a cada matriz de pesos empezando por la última y retrocediendo hasta llegar a la primera. Hacemos estos cálculos de forma vectorizada, procesando todos los casos a la vez, si bien cada matriz se calcula por separado en un bucle. La primera matriz que calculamos, que se corresponde a la última capa de la red, se obtiene aplicando las fórmulas siguientes:

$$\delta^{(L+1)} = a^{(L+1)} - y$$

$$D^{(L)} = \frac{(\delta^{(L+1)})^T a^{(L+1)} + \lambda \Theta_{\bullet, j>0}^{(L)}}{N}$$

El resto de matrices se calculan iterativamente con las fórmulas recursivas correspondientes, que para  $0 < i < L$ , donde  $L$  es el número de matrices de pesos que tiene la red, vienen dadas por:

$$\delta^{(i+1)} = (\delta^{(i+2)} \Theta(i+1)) \odot g'(z^{(i+1)})$$

$$D^{(i)} = \frac{(\delta_{\bullet, j>0}^{(i+1)})^T a^{(i)} + \lambda \Theta_{\bullet, j>0}^{(i)}}{N}$$

En estas fórmulas,  $N$  es el número de casos con el que se entrena la red,  $y$  es la respuesta esperada,  $a^{(l)}$  es la entrada de la capa  $l$  de la red,  $z$  es la salida de la capa  $l - 1$ ,  $\Theta^{(l)}$  es la matriz de pesos que pasa de la capa  $l$  a la  $l + 1$  y  $\lambda$  es el parámetro de regularización.

Una vez construida la función *gradiente*, utilizamos el código auxiliar proporcionado en el fichero *checkNNGradients.py* para comprobar que es correcta.

### 3. Resultados

Tras verificar que la función de coste y la de retro-propagación estaban implementadas correctamente, hemos procedido a entrenar la red neuronal con varios valores del parámetro de regularización  $\lambda$  y del número máximo de iteraciones para analizar el comportamiento de la red en cada caso. Para  $\lambda$ , hemos probado las potencias de 10 entre  $10^{-5}$  y  $10^2$ , y para el límite de iteraciones los múltiplos de 10 entre 10 y 100

$$\lambda \in \{10^i \mid i \in [-5, 2], i \in \mathbb{Z}\} \quad iter \in \{n \in [10, 100] \mid n \equiv 0 \pmod{10}\}$$

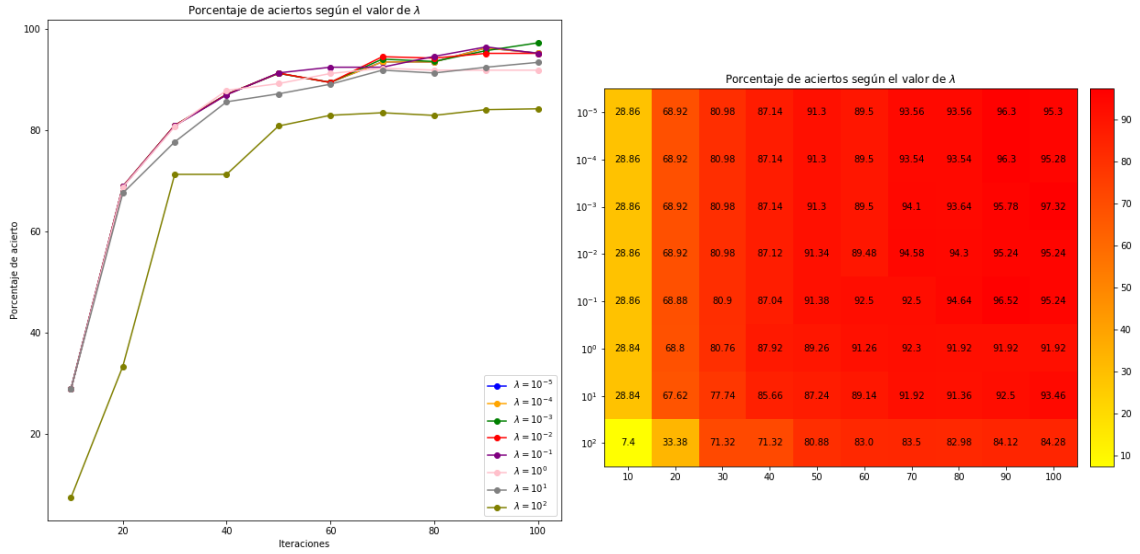


Figura 1: Estas gráficas muestran de dos formas distintas los mismos datos de tasas de acierto obtenidas tras entrenar nuestra red neuronal

Observamos que el porcentaje de aciertos de la red neuronal crece si permitimos un mayor número máximo de iteraciones. Esto se debe a que la función de minimización se va aproximando al mínimo progresivamente en cada iteración según le indica la dirección del gradiente. Si permitimos más iteraciones, nos aproximaremos más al mínimo de la función de coste y por tanto al mínimo error al predecir los datos de entrenamiento. Sin embargo, no es bueno minimizar excesivamente la función de coste puesto que, aunque permita acertar más con los datos de entrenamiento, puede perjudicar la predicción de datos nuevos si se produce sobreaprendizaje, haciendo que la red sólo prediga correctamente los datos de entrenamiento.

Asimismo, se obtienen más aciertos para valores más pequeños de  $\lambda$ . Para entender esto, recordamos que en el término de regularización que se añade a la función de coste,  $\lambda$  aparece multiplicando al sumatorio de los valores que aparecen en las matrices de pesos. Esto implica que, para minimizar la función de coste, mayores valores de  $\lambda$  fuerzan a elegir pesos más pequeños, mientras que  $\lambda$  más pequeños permiten mayor libertad al elegir los pesos. De este modo, con valores de  $\lambda$  pequeños la red tiene más flexibilidad para ajustarse a los datos de entrenamiento, y por tanto también existe la posibilidad del sobreaprendizaje si permitimos que la red se ajuste demasiado.

Por tanto, es posible que en los casos en los que nuestra red ha obtenido mejores tasas de acierto haya habido sobreaprendizaje, pero para confirmarlo necesitaríamos ver cómo se comporta la red con datos distintos de los utilizados para entrenarla.

## 4. Conclusiones

Gracias a esta práctica hemos comprobado que detrás del entrenamiento de una red neuronal hay muchos detalles técnicos y cálculos delicados que hay que realizar

con cuidado para que el resultado sea el deseado. Un fallo como sumar una columna de la matriz de pesos cuando no se debe puede pasar desapercibido al programador y provocar que los resultados de la red empeoren notablemente. En nuestro caso, al calcular el término de regularización de la función de coste incluimos las primeras columnas de las matrices de pesos, con lo que el resultado de la función de coste variaba ligeramente en el tercer decimal. Aunque al principio no dimos importancia a esta variación, provocó que el gradiente no se calculara bien y que el resultado no fuera igual que el obtenido por otro método implementado correctamente.

## 5. Anexo: Código

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

```
1 '''
2 Práctica 4
3
4 Rubén Ruperto Díaz y Rafael Herrera Troca
5 '''
6
7 import os
8 import numpy as np
9 import matplotlib.pyplot as plt
10 from matplotlib import cm
11 from mpl_toolkits.axes_grid1 import make_axes_locatable
12 from scipy.io import loadmat
13 from scipy.optimize import minimize
14
15
16 os.chdir("./resources")
17
18
19 # Función sigmoide
20 def sigmoide(z):
21     return 1 / (1 + np.exp(-z))
22
23 # Derivada de la función sigmoide
24 def diffSigmoide(a):
25     return a * (1 - a)
26
27 # Devuelve una matriz de pesos aleatorios con la dimensión dada
28 def randomWeights(l_in, l_out):
29     eps = np.sqrt(6)/np.sqrt(l_in + l_out)
30     rnd = np.random.random((l_out, l_in+1)) * (2*eps) - eps
31     return rnd
32
33 # Dada la entrada 'X' y los pesos 'theta' de una capa de una red
34 # neuronal, aplica los pesos y devuelve la salida de la capa
35 def applyLayer(X, theta):
36     thetaX = np.dot(X, theta.T)
37     return sigmoide(thetaX)
```

```

38
39 # Dada la entrada 'X' y el array de matrices de pesos 'theta',
40 # devuelve la entrada de cada capa y el resultado final devuelto
41 # por la red neuronal
42 def applyNet(X, theta):
43     lay = X.copy()
44     a = []
45     for i in range(len(theta)):
46         lay = np.hstack((np.array([np.ones(len(lay))]).T, lay))
47         a.append(lay.copy())
48         lay = applyLayer(lay, theta[i])
49
50     return lay,a
51
52 # Calcula la función de coste de una red neuronal para la
53 # salida esperada 'y', el resultado de la red 'h_theta', el array
54 # de matrices de pesos 'theta' y el término de regularización 'reg'
55 def coste(y, h_theta, theta, reg):
56     sumandos = -y * np.log(h_theta) - (1-y) * np.log(1-h_theta)
57     regul = 0
58     for i in range(len(theta)):
59         regul += np.sum(theta[i][:,1:]**2)
60     result = np.sum(sumandos) / len(y) + reg * regul / (2*len(y))
61     return result
62
63 # Calcula el gradiente de la función de coste haciendo
64 # retropropagación dada la salida esperada 'y', la entrada
65 # de cada capa 'a', la salida de la red 'h_theta', el array de
66 # matrices de pesos 'theta' y el término de regularización 'reg'
67 def gradiente(y, a, h_theta, theta, reg):
68     d = h_theta - y
69     delta = [np.dot(d.T, a[-1]) / len(y)]
70
71     for i in range(len(theta)-1,0,-1):
72         d = np.dot(d, theta[i]) * diffSigmoide(a[i])
73         d = d[:,1:]
74         delta.insert(0, np.dot(d.T, a[i-1]) / len(y))
75
76     for i in range(len(delta)):
77         delta[i][:,1:] += reg * theta[i][:,1:] / len(y)
78
79     return delta
80
81 # Calcula y devuelve el coste y el gradiente de una red neuronal
82 # dados todos los pesos en el array 'param_rn', las dimensiones
83 # de cada capa en 'capas', los datos de entrada 'X', la salida
84 # esperada 'y' y el término de regularización 'reg'
85 def backprop(params_rn, capas, X, y, reg):
86     # Convertimos el vector de todos los pesos en las distintas
87     # matrices
88     theta = [np.reshape(params_rn[:capas[1]*(capas[0]+1)],(capas
89 [1],capas[0]+1))]
89     gastados = capas[1]*(capas[0]+1)
90     for i in range(len(capas)-2):

```

```

91     theta.append(np.reshape(params_rn[gastados:gastados+capas[i
+2]*(capas[i+1]+1)],(capas[i+2],capas[i+1]+1)))
92     gastados += capas[i+2]*(capas[i+1]+1)
93
94     # Calculamos el vector de salida esperada para la red neuronal
95     Y = np.zeros((len(y), capas[-1]))
96     for i in range(len(Y)):
97         Y[i,y[i]-1] = 1
98
99     # Aplicamos la red neuronal
100    h_theta,a = applyNet(X, theta)
101
102    cost = coste(Y, h_theta, theta, reg)
103    grad = gradiente(Y, a, h_theta, theta, reg)
104
105    g = np.array([])
106    for i in range(len(grad)):
107        g = np.concatenate((g, grad[i].ravel()))
108
109    return cost, g
110
111 # Calcula el porcentaje de acierto obtenido con la respuesta dada
112 # en 'X' y las etiquetas correctas en 'Y'
113 def acierto(X, Y):
114     resultados = X.argmax(axis=1) + 1
115     return 100 * np.count_nonzero(resultados == Y.ravel()) / len(Y)
116
117
118 ## Parte 1: Función de coste y gradiente
119
120 # Leemos los datos
121 data = loadmat('ex4data1.mat')
122 y = data['y'].ravel()
123 X = data['X']
124
125 # Leemos las matrices de pesos
126 weights = loadmat('ex4weights.mat')
127 theta1, theta2 = weights['Theta1'], weights['Theta2']
128
129 # Calculamos el coste y el gradiente
130 theta = np.concatenate((theta1.ravel(), theta2.ravel()))
131 res = backprop(theta, (400,25,10), X, y, 1)
132
133
134 ## Parte 2: Entrenamos la red neuronal
135
136 # Creamos unas matrices de pesos iniciales de forma aleatoria
137 theta01 = randomWeights(400, 25)
138 theta02 = randomWeights(25, 10)
139 theta0 = np.concatenate((theta01.ravel(), theta02.ravel()))
140
141 # Entrenamos la red neuronal con distintos valores para el
142 # parámetro de regularización y representamos su porcentaje
143 # de acierto
144 reg = range(-5, 3)

```



```

145 itera = range(10, 110, 10)
146
147 pAc = np.zeros((len(reg), len(itera)))
148 for i in range(len(reg)):
149     for j in range(len(itera)):
150         theta = minimize(fun=backprop, x0=theta0, args=((400,25,10)
151             , X, y, 10**reg[i]), method='TNC', jac=True, options={'maxiter':
152                 itera[j]})['x']
153
154         theta1 = np.reshape(theta[:25*401], (25,401))
155         theta2 = np.reshape(theta[25*401:], (10,26))
156         res = applyNet(X, (theta1, theta2))[0]
157
158         pAc[i][j] = acierto(res, y)
159
160 colores = ['blue', 'orange', 'green', 'red', 'purple', 'pink', '
161     gray', 'olive', 'cyan']
162 plt.figure(figsize=(10,10))
163 for i in range(len(pAc)):
164     plt.plot(range(10, 110, 10), pAc[i], color=colores[i], marker='
165         o', label=(r"$\lambda = 10^{" + str(reg[i]) + "}"$"))
166 plt.title(r"Porcentaje de aciertos según el valor de $\lambda$")
167 plt.xlabel("Iteraciones")
168 plt.ylabel("Porcentaje de acierto")
169 plt.legend(loc="lower right")
170 plt.savefig("LineasLambda.png")
171 plt.show()
172
173 xLabels = [str(it) for it in itera]
174 yLabels = [r'$10^{'+ str(r) + '}$' for r in reg]
175 plt.figure(figsize=(10,10))
176 plt.title(r"Porcentaje de aciertos según el valor de $\lambda$")
177 fig = plt.subplot()
178 im = fig.imshow(pAc, cmap=cm.autumn_r)
179 cax = make_axes_locatable(fig).append_axes("right", size="5%", pad
180     =0.2)
181 plt.colorbar(im, cax=cax)
182 for i in range(len(xLabels)):
183     for j in range(len(yLabels)):
184         text = fig.text(i, j, pAc[j][i], ha="center", va="center",
185             color="k")
186 fig.set_xticks(np.arange(len(xLabels)))
187 fig.set_yticks(np.arange(len(yLabels)))
188 fig.set_xticklabels(xLabels)
189 fig.set_yticklabels(yLabels)
190 plt.savefig("CuadrosLambda.png")
191 plt.show()

```