# Neural Networks

## Matt R

## February 21, 2022

# Contents

# 1 Logistic Regression

We begin with a review of binary classification and logistic regression. To this end, suppose we have we have training examples $x \in \mathbb{R}^{m \times n}$ with binary labels $y \in \{0,1\}^{1 \times n}$. We desire to train a model which yields an output $a$ which represents

$$a = \mathbb{P}(y = 1 | x).$$

To this end, let $\sigma : \mathbb{R} \to (0,1)$ denote the sigmoid function, i.e.,

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

and let $w \in \mathbb{R}^m$, $b \in \mathbb{R}$, and let

$$a = \sigma(w^T x + b).$$

To analyze the accuracy of model, we need a way to compare $y$ and $a$, and ideally this functional comparison can be optimized with respect to $(w, b)$ in such a way to minimize the error. To this end, we note that

$$\mathbb{P}(y|x) = a^y (1 - a)^{1-y},$$

or rather

$$\mathbb{P}(y = 1|x) = a, \qquad \mathbb{P}(y = 0|x) = 1 - a,$$

so $\mathbb{P}(y|x)$ represents the corrected probability. Now since we want

$$a \approx 1 \quad \text{when } y = 1,$$

and

$$a \approx 0 \quad \text{when } y = 0,$$

and $0 \leq a \leq 1$, any error using differences won't be refined enough to analyze when tuning the model. Moreover, since introducing the sigmoid function, our usual mean-squared-error function won't be convex. This leads us to apply the log function, which when restricted to $(0,1)$ is a bijective mapping of $(0,1) \to (-\infty, 0)$. This leads us to define our log-loss function

$$\begin{aligned}
\mathbb{L}(a, y) &= -\log(\mathbb{P}(y|x)) \\
&= -\log\left(a^y (1 - a)^{1-y}\right) \\
&= -\left[y \log(a) + (1 - y) \log(1 - a)\right],
\end{aligned}$$

and finally, since we wish to analyze how our model performs on the entire training set, we need to average our log-loss functions to obtain our cost function $\mathbb{J}$ defined by

$$\begin{aligned}
\mathbb{J}(w, b) &= \frac{1}{n} \sum_{j=1}^{n} \mathbb{L}(a_j, y_j) \\
&= -\frac{1}{n} \sum_{j=1}^{n} \left[ y_j \log(a_j) + (1 - y_j) \log(1 - a_j) \right] \\
&= -\frac{1}{n} \sum_{j=1}^{n} \left[ y_j \log(\sigma(w^T x_j + b)) + (1 - y_j) \log(1 - \sigma(w^T x_j + b)) \right].
\end{aligned}$$

## 1.1  The Gradient

To compute the gradient of our cost function $\mathbb{J}$, we first write $\mathbb{J}$ as a sum of compositions as follows: We have the log-loss function considered as a map $\mathbb{L} : (0, 1) \times \mathbb{R} \to \mathbb{R}$,

$$\mathbb{L}(a, y) = -\left[ y \log(a) + (1 - y) \log(1 - a) \right],$$

we have the sigmoid function $\sigma : \mathbb{R} \to (0, 1)$ with $\sigma(z) = a$ and $\sigma'(z) = a(1 - a)$, and we have the collection of affine-functionals $\phi_x : \mathbb{R}^m \times \mathbb{R} \to \mathbb{R}$ given by

$$\phi_x(w, b) = w^T x + b,$$

for which we fix an arbitrary $x \in \mathbb{R}^m$ and write $\phi = \phi_x$, and set $z = \phi(w, b)$. Finally, we introduce the auxiliary function $\mathcal{L} : \mathbb{R}^m \times \mathbb{R} \to \mathbb{R}$ given by

$$\mathcal{L}(w, b) = \mathbb{L}(\sigma(\phi(w, b)), y).$$

Then by the chain rule, we have that

$$\begin{aligned}
d\mathcal{L} &= d_a \mathbb{L}(a, y) \circ d\sigma(z) \circ d_w \phi(w, b) \\
&= \left[ -\frac{y}{a} + \frac{1 - y}{1 - a} \right] \cdot a(1 - a) \cdot \begin{bmatrix} x^T & 1 \end{bmatrix} \\
&= \left[ -y(1 - a) + a(1 - y) \right] \cdot \begin{bmatrix} x^T & 1 \end{bmatrix} \\
&= (a - y) \begin{bmatrix} x^T & 1 \end{bmatrix}
\end{aligned}$$

Composition turns into matrix multiplication in the tangent space.

Moreover, since in Euclidean space, we have that $\nabla f = (df)^T$, and hence that

$$\nabla \mathcal{L}(w, b) = (a - y) \begin{bmatrix} x \\ 1 \end{bmatrix},$$

or rather

$$\partial_w \mathbb{L}(a, y) = (a - y)x, \qquad \partial_b \mathbb{L}(a, y) = a - y.$$

Finally, since our cost function $\mathbb{J}$ is the sum-log-loss, we have by linearity that

$$\partial_w \mathbb{J}(w, b) = \frac{1}{n} \sum_{j=1}^{n} (a_j - y_j)x_j$$

$$= \frac{1}{n}((a - y) \cdot x^T)^T$$

$$= \frac{1}{n} x \cdot (a - y)^T$$

and

$$\partial_b \mathbb{J}(w, b) = \frac{1}{n} \sum_{j=1}^{n} (a_j - y_j).$$

### 1.1.1 Vectorization in Python

Here we include the general code to train a model using logistic regression without regularization and without tuning on a cross-validation set.

```python
import copy

import numpy as np

def sigmoid(z):
    """
    Parameters
    ----------
    z : array_like

    Returns
    -------
    sigma : array_like
    """

    sigma = (1 / (1 + np.exp(-z)))
    return sigma

```

```python
19  def cost_function(x, y, w, b):
20      """
21      Parameters
22      ----------
23      x : array_like
24          x.shape = (m, n) with m-features and n-examples
25      y : array_like
26          y.shape = (1, n)
27      w : array_like
28          w.shape = (m, 1)
29      b : float
30
31      Returns
32      -------
33      J : float
34          The value of the cost function evaluated at (w, b)
35      dw : array_like
36          dw.shape = w.shape = (m, 1)
37          The gradient of J with respect to w
38      db : float
39          The partial derivative of J with respect to b
40      """
41
42      # Auxiliary assignments
43      m, n = x.shape
44      z = w.T @ x + b
45      assert z.size == n
46      a = sigmoid(z).reshape(1, n)
47      dz = a - y
48
49      # Compute cost J
50      J = (-1 / n) * (np.log(a) @ y.T + np.log(1 - a) @ (1 - y).T)
51
52      # Compute dw and db
53      dw = (x @ dz.T) / m
54      assert dw.shape == w.shape
55      db = np.sum(dz) / m
56
57      return J, dw, db
58
59  def grad_descent(x, y, w, b, alpha=0.001, num_iters=2000, print_cost=False):
60      """
61      Parameters
62      ----------
63      x, y, w, b : See cost_function above for specifics.
64          w and b are chosen to initialize the descent (likely all components 0)
65      alpha : float
```

```
66          The learning rate of gradient descent
67      num_iters : int
68          The number of times we wish to perform gradient descent
69
70      Returns
71      -------
72      costs : List[float]
73          For each iteration we record the cost-values associated to (w, b)
74      params : Dict[w : array_like, b : float]
75          w : array_like
76              Optimized weight parameter w after iterating through grad descent
77          b : float
78              Optimized bias parameter b after iterating through grad descent
79      grads : Dict[dw : array_like, db : float]
80          dw : array_like
81              The optimized gradient with repsect to w
82          db : float
83              The optimized derivative with respect to b
84      """
85
86      costs = []
87      w = copy.deepcopy(w)
88      b = copy.deepcopy(b)
89      for i in range(num_iters):
90          J, dw, db = cost_function(x, y, w, b)
91          w = w - alpha * dw
92          b = b - alpha * db
93
94          if i % 100 == 0:
95              costs.append(J)
96              if print_cost:
97                  idx = int(i / 100) - 1
98                  print(f'Cost after iteration {i}: {costs[idx]}')
99
100     params = {'w' : w, 'b' : b}
101     grads = {'dw' : dw, 'db' : db}
102
103     return costs, params, grads
104
105 def predict(w, b, x):
106     """
107     Parameters
108     ----------
109     w : array_like
110         w.shape = (m, 1)
111     b : float
112     x : array_like
```

```
113        x.shape = (m, n)

115    Returns
116    -------
117    y_predict : array_like
118        y_pred.shape = (1, n)
119        An array containing the prediction of our model applied to training
120        data x, i.e., y_pred = 1 or y_pred = 0.
121    """

123    m, n = x.shape
124    # Get probability array
125    a = sigmoid(w.T @ x + b)
126    # Get boolean array with False given by a < 0.5
127    pseudo_predict = ~(a < 0.5)
128    # Convert to binary to get predictions
129    y_predict = pseudo_predict.astype(int)

131    return y_predict

133 def model(x_train, y_train, x_test, y_test, alpha=0.001, num_iters=2000, accuracy=Tr
134    """
135    Parameters:
136    -----------
137    x_train, y_train, x_test, y_test : array_like
138        x_train.shape = (m, n_train)
139        y_train.shape = (1, n_train)
140        x_test.shape = (m, n_test)
141        y_test.shape = (1, n_test)
142    alpha : float
143        The learning rate for gradient descent
144    num_iters : int
145        The number of times we wish to perform gradient descent
146    accuracy : Boolean
147        Use True to print the accuracy of the model

149    Returns:
150    d : Dict
151        d['costs'] : array_like
152            The costs evaluated every 100 iterations
153        d['y_train_preds'] : array_like
154            Predicted values on the training set
155        d['y_test_preds'] : array_like
156            Predicted values on the test set
157        d['w'] : array_like
158            Optimized parameter w
159        d['b'] : float
```

```
160              Optimized parameter b
161        d['learning_rate'] : float
162              The learning rate alpha
163        d['num_iters'] : int
164              The number of iterations with which gradient descent was performed
165
166     """
167
168     m = x_train.shape[0]
169     # initialize parameters
170     w = np.zeros((m, 1))
171     b = 0.0
172     # optimize parameters
173     costs, params, grads = grad_descent(x_train, y_train, w, b, alpha, num_iters)
174     w = params['w']
175     b = params['b']
176     # record predictions
177     y_train_preds = predict(w, b, x_train)
178     y_test_preds = predict(w, b, x_test)
179     # group results into dictionary for return
180     d = {'costs' : costs,
181         'y_train_preds' : y_train_preds,
182         'y_test_preds' : y_test_preds,
183         'w' : w,
184         'b' : b,
185         'learning_rate' : alpha,
186         'num_iters' : num_iters}
187
188     if accuracy:
189         train_acc = 100 - np.mean(np.abs(y_train_preds - y_train)) * 100
190         test_acc = 100 - np.mean(np.abs(y_test_preds - y_test)) * 100
191         print(f'Training_Accuracy:_{train_acc}%')
192         print(f'Test_Accuracy:_{test_acc}%')
193
194
195     return d
```