

Neural Networks

Matt R

June 28, 2022

Contents

I	Neural Networks and Deep Learning	2
1	Logistic Regression	3
1.1	The Gradient	4
1.2	Implementation in Python via <code>numpy</code>	7
1.3	Implementation in Python via <code>sklearn</code>	11
2	Neural Networks: A Single Hidden Layer	13
2.1	Activation Functions	15
2.1.1	The Sigmoid Function	15
2.1.2	The Hyperbolic Tangent Function	16
2.1.3	The Rectified Linear Unit Function	16
2.1.4	The Softmax Function	17
2.2	Backward Propagation	18
3	Deep Neural Networks	25
3.1	Backward Propagation	25
3.2	Implementation in Python via <code>numpy</code>	26
3.3	Implementation in Python via <code>tensorflow</code>	30
	References	33

Part I

Neural Networks and Deep Learning

1 Logistic Regression

We begin with a review of binary classification and logistic regression. To this end, suppose we have training examples $x \in \mathbb{R}^{n \times N}$ 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 $\sigma : \mathbb{R} \rightarrow (0, 1)$ denote the sigmoid function, i.e.,

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

and let $w \in \mathbb{R}^{1 \times n}$, $b \in \mathbb{R}$, and let

$$a = \sigma(wx + b).$$

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 an error. To this end, we note that

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

or rather

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

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 $(0, 1) \rightarrow (-\infty, 0)$. This leads us to define our log-loss function

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

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

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

1.1 The Gradient

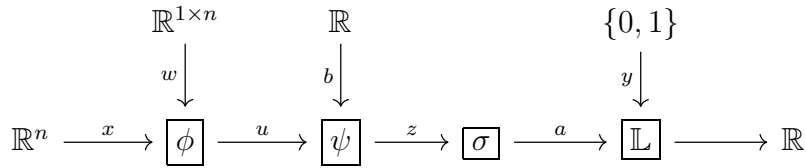
We wish to compute the gradient of our cost function \mathbb{J} with respect to our trainable parameters, $w \in \mathbb{R}^{1 \times n}$ and $b \in \mathbb{R}$. To this end, we define the functions

$$\phi : \mathbb{R}^{1 \times n} \times \mathbb{R}^n \rightarrow \mathbb{R}, \quad \phi(w, x) = wx,$$

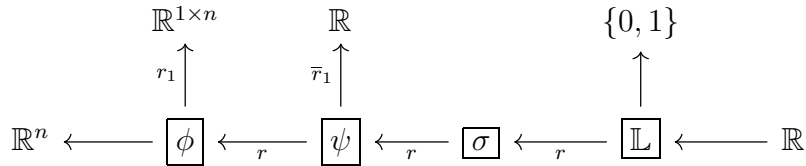
and

$$\psi : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}, \quad \psi(b, u) = u + b.$$

Then our logistic regression model for a single example follows the following network layout:



Let's now analyze our reverse differentials for this type of composition:



1.

$$\phi : \mathbb{R}^{1 \times n} \times \mathbb{R}^n \rightarrow \mathbb{R}, \quad u := \phi(w, x) = wx.$$

Then for any $(w, x) \in \mathbb{R}^{1 \times n} \times \mathbb{R}^n$ and any $\eta \in T_w \mathbb{R}^{1 \times n}$, we have that

$$\begin{aligned} d_1 \phi_{(w, x)}(\eta) &= \eta x \\ &= R_x(\eta), \end{aligned}$$

where R_x is the right-multiplication operator. It then follows that for any $\zeta \in T_u \mathbb{R}$, that

$$\begin{aligned} \langle r_1 \phi_{(w, x)}(\zeta), \eta \rangle_{\mathbb{R}^{1 \times n}} &= \langle \zeta, d_1 \phi_{(w, x)}(\eta) \rangle_{\mathbb{R}} \\ &= \langle \zeta, R_x(\eta) \rangle_{\mathbb{R}} \\ &= \langle R_{x^T}(\zeta), \eta \rangle_{\mathbb{R}^{1 \times n}}, \end{aligned}$$

and hence that

$$r_1 \phi_{(w, x)} = R_{x^T}.$$

2.

$$\psi : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}, \quad z := \psi(b, u) = u + b.$$

Then for any $(b, u) \in \mathbb{R} \times \mathbb{R}$ and any $\xi \in T_u \mathbb{R}$, we have that

$$d\psi_{(b, u)}(\xi) = \mathbf{1}_{\mathbb{R}}(\xi),$$

and similarly for any $\eta \in T_b \mathbb{R}$, we have that

$$\bar{d}_1 \psi_{(b, u)}(\eta) = \mathbf{1}_{\mathbb{R}}(\eta).$$

We then immediately have that

$$r\psi_{(b, u)} = \mathbf{1}_{\mathbb{R}},$$

and

$$\bar{r}_1 \psi_{(b, u)} = \mathbf{1}_{\mathbb{R}}.$$

3.

$$\sigma : \mathbb{R} \rightarrow \mathbb{R}, \quad a := \sigma(z) = \frac{1}{1 + e^{-z}}.$$

Then

$$\begin{aligned}
r\sigma_z &= \frac{e^{-z}}{(1 + e^{-z})^2} \\
&= \frac{1}{1 + e^{-z}} \frac{e^{-z}}{1 + e^{-z}} \\
&= \sigma(z) \frac{1 + e^{-z} - 1}{1 + e^{-z}} \\
&= \sigma(z)(1 - \sigma(z)) \\
&= a(1 - a).
\end{aligned}$$

4.

$$\mathbb{L} : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}, \quad \mathbb{L}(a, y) = -[y \log(a) + (1 - y) \log(1 - a)].$$

Then

$$r\mathbb{L}_{(a,y)} = -\frac{y}{a} + \frac{1-y}{1-a}$$

We now compute the gradients with respect to w and b . To this end,

$$\begin{aligned}
\frac{\partial \mathbb{J}}{\partial w} &= \frac{1}{N} \sum_{j=1}^N r_1 \phi_{w, x_j} \circ r\psi_{(b, u_j)} \circ r\sigma_{z_j} \circ r\mathbb{L}_{(a_j, y_j)} \\
&= \frac{1}{N} \sum_{j=1}^N R_{x_j^T} \circ \left[-\frac{y_j}{a_j} + \frac{1-y_j}{1-a_j} \right] \cdot (a_j(1-a_j)) \\
&= \frac{1}{N} \sum_{j=1}^N (a_j - y_j) x_j^T \\
&= \frac{1}{N} (a - y) x^T,
\end{aligned}$$

and

$$\begin{aligned}
\frac{\partial \mathbb{J}}{\partial b} &= \frac{1}{N} \sum_{j=1}^N \bar{r}_1 \psi_{b, u_j} \circ r\sigma_{z_j} \circ r\mathbb{L}_{(a_j, y_j)} \\
&= \frac{1}{N} \sum_{j=1}^N (a_j - y_j)
\end{aligned}$$

1.2 Implementation in Python via numpy

Here we include the general method of coding a logistic regression model with L^2 -regularization via the classical numpy library.

```
1  #! python3
2
3  import numpy as np
4
5  from mllib.utils import apply_activation
6
7  class LinearParameters():
8      def __init__(self, dims, bias=True, seed=1):
9          """
10         Parameters:
11         -----
12         dims : tuple(int, int)
13         bias : Boolean
14             Default : True
15         seed : int
16             Default : 1
17
18         Returns:
19         -----
20         None
21         """
22         np.random.seed(seed)
23         self.dims = dims
24         self.bias = bias
25         self.w = np.random.randn(*dims) * 0.01
26         if bias:
27             self.b = np.zeros((dims[0], 1))
28
29     def forward(self, x):
30         """
31         Parameters:
32         -----
33         x : array_like
34
35         Returns:
36         -----
37         z : array_like
38         """
39         z = np.einsum('ij,jk', self.w, x)
40         if self.bias:
41             z += self.b
42
```

```

43         return z
44
45     def backward(self, dz, x):
46         """
47         Parameters:
48         -----
49         dz : array_like
50         x : array_like
51
52         Returns:
53         -----
54         None
55         """
56         if self.bias:
57             self.db = np.sum(dz, axis=1, keepdims=True)
58             assert (self.db.shape == self.b.shape)
59
60             self.dw = np.einsum('ij,kj', dz, x)
61             assert (self.dw.shape == self.w.shape)
62
63     def update(self, learning_rate=0.01):
64         """
65         Parameters:
66         -----
67         learning_rate : float
68             Default : 0.01
69
70         Returns:
71         -----
72         None
73         """
74         w = self.w - learning_rate * self.dw
75         self.w = w
76
77         if self.bias:
78             b = self.b - learning_rate * self.db
79             self.b = b
80
81 class LogisticRegression():
82     def __init__(self, lp_reg):
83         """
84         Parameters:
85         lp_reg : int
86             2 : L_2 Regularization is imposed
87             1 : L_1 Regularization is imposed
88             0 : No regulariation is imposed
89

```



```

90         Returns:
91         -----
92         None
93         """
94         self.lp_reg = lp_reg
95
96     def predict(self, params, x):
97         """
98         Parameters:
99         -----
100         params : class[LinearParameters]
101         x : array_like
102
103         Returns:
104         -----
105         a : array_like
106         dg : array_like
107         """
108         z = params.forward(x)
109         a, dg = apply_activation(z, 'sigmoid')
110         return a, dg
111
112     def cost_function(self, params, x, y, lambda_=0.01, eps=1e-8):
113         """
114         Parameters:
115         -----
116         params : class[LinearParameters]
117         x : array_like
118         y : array_like
119         lambda_ : float
120             Default : 0.01
121         eps : float
122             Default : 1e-8
123
124         Returns:
125         -----
126         cost : float
127         """
128         n = y.shape[1]
129
130         R = np.sum(np.abs(params.w) ** self.lp_reg)
131         R *= (lambda_ / (2 * n))
132
133         a, _ = self.predict(params, x)
134         a = np.clip(a, eps, 1 - eps)
135
136         J = (-1 / n) * (np.sum(y * np.log(a) + (1 - y) * np.log(1 - a)))

```

```

137
138         cost = float(np.squeeze(J + R))
139
140     return cost
141
142     def fit(self, x, y, learning_rate=0.1, lambda_=0.01, seed=1, num_iters=10000):
143         """
144         Parameters:
145         -----
146         x : array_like
147         y : array_like
148         learning_rate : float
149             Default : 0.1
150         lambda_ : float
151             Default : 0.0
152         num_iters : int
153             Default : 10000
154
155         Returns:
156         -----
157         costs : List[floats]
158         params : class[Parameters]
159         """
160         dims = (y.shape[0], x.shape[0])
161         n = x.shape[1]
162         params = LinearParameters(dims, True, seed)
163
164         if self.lp_reg == 0:
165             lambda_ = 0.0
166
167         costs = []
168         for i in range(num_iters):
169             a, _ = self.predict(params, x)
170             cost = self.cost_function(params, x, y, lambda_)
171             costs.append(cost)
172             dz = (a - y) / n
173             params.backward(dz, x)
174             params.update(learning_rate)
175
176             if i % 1000 == 0:
177                 print(f'Cost_after_iteration_{i}:_{cost}')
178
179         return params
180
181     def evaluate(self, params, x):
182         """
183         Parameters:

```

```

184         -----
185         params : class[Parameters]
186         x : array_like
187
188         Returns:
189         -----
190         y_hat : array_like
191         """
192         a, _ = self.predict(params, x)
193         y_hat = (~(a < 0.5)).astype(int)
194
195         return y_hat
196
197     def accuracy(self, params, x, y):
198         """
199         Parameters:
200         -----
201         params : class[Parameters]
202         x : array_like
203         y : array_like
204
205         Returns:
206         -----
207         accuracy : float
208         """
209         y_hat = self.evaluate(params, x)
210
211         accuracy = np.sum(y_hat == y) / y.shape[1]
212
213         return accuracy

```

1.3 Implementation in Python via sklearn

Here we include the general method of coding a logistic regression model via scikit-learn's modeling library.

```

1  #! python3
2
3  import pandas as pd
4  import numpy as np
5  from sklearn.model_selection import train_test_split
6  from sklearn.linear_model import LogisticRegression
7
8  def main(csv):
9      df = pd.read_csv(csv)
10     dataset = df.values

```

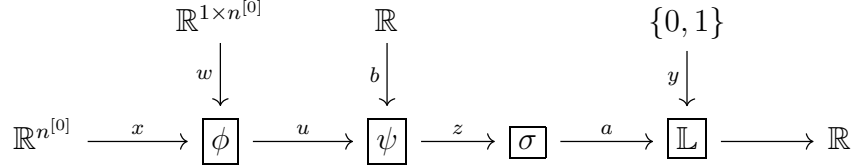
```

11     x = dataset[:, :10]
12     y = dataset[:, 10]
13
14     x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2)
15     mu = np.mean(x, axis=0, keepdims=True)
16     var = np.var(x, axis=0, keepdims=True)
17     x_train = (x_train - mu) / np.sqrt(var)
18     x_test = (x_test - mu) / np.sqrt(var)
19
20     log_reg = LogisticRegression()
21     log_reg.fit(x_train, y_train)
22     train_acc = log_reg.score(x_train, y_train)
23     print(f'The accuracy on the training set: {train_acc}.')
24     test_acc = log_reg.score(x_test, y_test)
25     print(f'The accuracy on the test set: {test_acc}.')

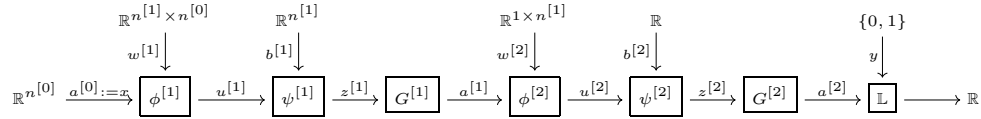
```

2 Neural Networks: A Single Hidden Layer

Suppose we wish to consider the binary classification problem given the training set (x, y) with $x \in \mathbb{R}^{n^{[0]} \times N}$ and $y \in \{0, 1\}^{1 \times N}$. Usually with logistic regression we have the following type of structure:



Such a structure will be called a *network*, and the a is known as the *activation node*. Logistic regression can be too simplistic of a model for many situations, e.g., if the dataset isn't linearly separable (i.e., there doesn't exist some well-defined decision boundary built from a linear-surface), then logistic regression won't give a high-accuracy model. To modify this model to handle more complex situations, we introduce a new "hidden layer" of nodes with their own (possibly different) activation functions. That is, we consider a network of the following form:



In the above diagram, we use $\cdot^{[0]}$ to denote everything in layer-0, i.e., the input layer; we use $\cdot^{[1]}$ to denote everything in layer-1, i.e., the hidden layer; and we use $\cdot^{[2]}$ to denote everything in layer-2, i.e., the output layer. Moreover, we have the functions (where we suppress the layer-notation)

- $\phi : \mathbb{R}^{n \times m} \times \mathbb{R}^m \rightarrow \mathbb{R}^n, \quad u := \phi(w, a) = wa,$
- $\psi : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}^n, \quad z := \psi(b, u) = u + b,$
- $G : \mathbb{R}^n \rightarrow \mathbb{R}^n, \quad a := G(z),$

where G is the broadcasting of something activation function $g : \mathbb{R} \rightarrow \mathbb{R}$.

Definition 2.1. Suppose $g : \mathbb{R} \rightarrow \mathbb{R}$ is any function. Then we say $G : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is the **broadcast** of g from \mathbb{R} to \mathbb{R}^n if

$$\begin{aligned} G(v) &= G(v^i e_i) \\ &= g(v^i) e_i, \end{aligned}$$

where $v \in \mathbb{R}^n$ and $\{e_i : 1 \leq i \leq n\}$ is the standard basis for \mathbb{R}^n . In practice, we will sometimes write $g = G$ for a broadcasted function, and let the context determine the meaning of g .

castingDifferential

Lemma 2.2. Suppose $g : \mathbb{R} \rightarrow \mathbb{R}$ is any smooth function and $G : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is the broadcasting of g from \mathbb{R} to \mathbb{R}^n . Then the differential $dG_z : T_z \mathbb{R}^n \rightarrow T_{G(z)} \mathbb{R}^n$ is given by

$$dG_z(\xi) = [g'(z^i)] \odot [\xi^i],$$

where \odot is the Hadamard product (also know as component-wise multiplication), and has matrix-representation in $\mathbb{R}^{m \times m}$ given by

$$[dG_z]_j^i = \delta_j^i g'(z^i).$$

We use the notation

$$G'(z) := [g'(z^i)] \in \mathbb{R}^n,$$

and thus may write

$$dG_z(v) = G'(z) \odot \xi.$$

Furthermore, we have that for $\zeta \in T_{G(z)} \mathbb{R}^n$,

$$rG_z(\zeta) = G'(z) \odot \zeta.$$

Proof: We calculate

$$\begin{aligned} dG_z(\xi) &= \left. \frac{d}{dt} \right|_{t=0} G(z + t\xi) \\ &= \left. \frac{d}{dt} \right|_{t=0} (g(z^i + t\xi^i)) \\ &= (g'(z^i) \xi^i) \\ &= [g'(z^i)] \odot [\xi^i], \end{aligned}$$

and letting e_1, \dots, e_m denote the usual basis for $T_z \mathbb{R}^m$ (identified with \mathbb{R}^m), we see that

$$\begin{aligned} dG_z(e_j) &= [g'(z^i)] \odot e_j \\ &= g'(z^j) e_j, \end{aligned}$$

from which conclude that dG_z is diagonal with (j, j) -th entry $g'(z^j)$ as desired.

Furthermore, for $\zeta \in T_{G(z)}\mathbb{R}^n$, we have that

$$\begin{aligned}\langle rG_z(\zeta), \xi \rangle_{\mathbb{R}^n} &= \langle \zeta, dG_z(\xi) \rangle_{\mathbb{R}^n} \\ &= \langle \zeta, G'(z) \odot \xi \rangle_{\mathbb{R}^n} \\ &= \langle G'(z) \odot \zeta, \xi \rangle_{\mathbb{R}^n},\end{aligned}$$

and the result follows. \square

Returning to our network, we see call the full composition of network functions resulting in $a^{[2]}$, the *forward propagation*. That is, given an example $x \in \mathbb{R}^{n^{[0]}}$, we have that

$$a^{[2]} = G^{[2]}(\psi^{[2]}(b^{[2]}, \phi^{[2]}(w^{[2]}, G^{[1]}(\psi^{[1]}(b^{[1]}, \phi^{[1]}(w^{[1]}, x)))))).$$

2.1 Activation Functions

There are mainly only a handful of activating functions we consider for our non-linearity conditions (but many more built from these that follow).

2.1.1 The Sigmoid Function

We have the sigmoid function $\sigma(z)$ given by

$$\sigma : \mathbb{R} \rightarrow (0, 1), \quad \sigma(z) = \frac{1}{1 + e^{-z}}.$$

We note that since

$$\begin{aligned}1 - \sigma(z) &= 1 - \frac{1}{1 + e^{-z}} \\ &= \frac{e^{-z}}{1 + e^{-z}}\end{aligned}$$

$$\begin{aligned}\sigma'(z) &= \frac{e^{-z}}{(1 + e^{-z})^2} \\ &= \frac{1}{1 + e^{-z}} \cdot \frac{e^{-z}}{1 + e^{-z}} \\ &= \sigma(z)(1 - \sigma(z))\end{aligned}$$

2.1.2 The Hyperbolic Tangent Function

We have the hyperbolic tangent function $\tanh(z)$ given by

$$\tanh : \mathbb{R} \rightarrow (-1, 1), \quad \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}.$$

We then calculate

$$\begin{aligned} \tanh'(z) &= \frac{(e^z + e^{-z})(e^z + e^{-z}) - (e^z - e^{-z})(e^z - e^{-z})}{(e^z + e^{-z})^2} \\ &= \frac{(e^z + e^{-z})^2}{(e^z + e^{-z})^2} - \frac{(e^z - e^{-z})^2}{(e^z + e^{-z})^2} \\ &= 1 - \tanh^2(z). \end{aligned}$$

Furthermore, we note that

$$\frac{1}{2} \left(\tanh\left(\frac{z}{2}\right) + 1 \right) = \sigma(z).$$

Indeed,

$$\begin{aligned} 1 + \tanh \frac{z}{2} &= 1 + \frac{e^{\frac{z}{2}} - e^{-\frac{z}{2}}}{e^{\frac{z}{2}} + e^{-\frac{z}{2}}} \\ &= \frac{e^{\frac{z}{2}} + e^{-\frac{z}{2}} + e^{\frac{z}{2}} - e^{-\frac{z}{2}}}{e^{\frac{z}{2}} + e^{-\frac{z}{2}}} \\ &= 2 \frac{e^{\frac{z}{2}}}{e^{\frac{z}{2}} + e^{-\frac{z}{2}}} \\ &= 2 \frac{1}{1 + e^{-z}} \\ &= 2\sigma(z), \end{aligned}$$

as desired.

2.1.3 The Rectified Linear Unit Function

We have the leaky-ReLU function $\text{ReLU}(z; \beta)$ given by

$$\text{ReLU} : \mathbb{R} \rightarrow \mathbb{R}, \quad \text{ReLU}(z; \beta) = \max\{\beta z, z\},$$

for some $\beta > 0$ (typically chosen very small).

We have the rectified linear unit function $\text{ReLU}(z)$ given by setting $\beta = 0$ in the leaky-ReLu function, i.e.,

$$\text{ReLU} : \mathbb{R} \rightarrow [0, \infty), \quad \text{ReLU}(z) = \text{ReLU}(z; \beta = 0) = \max\{0, z\}.$$

We then calculate

$$\begin{aligned} \text{ReLU}'(z; \beta) &= \begin{cases} \beta & z < 0 \\ 1 & z \geq 0 \end{cases} \\ &= \beta \chi_{(-\infty, 0)}(z) + \chi_{[0, \infty)}(z), \end{aligned}$$

where

$$\chi_A(z) = \begin{cases} 1 & z \in A \\ 0 & z \notin A \end{cases},$$

is the indicator function.

2.1.4 The Softmax Function

We finally have the softmax function $\text{softmax}(z)$ given by

$$\text{softmax} : \mathbb{R}^m \rightarrow \mathbb{R}^m, \quad \text{softmax}(z) = \frac{1}{\sum_{j=1}^m e^{z^j}} \begin{pmatrix} e^{z^1} \\ e^{z^2} \\ \vdots \\ e^{z^m} \end{pmatrix},$$

which we typically use this function on the outer-layer to obtain a probability distribution over our predicted labels when dealing with multi-class regression. Let

$$S^i = x^i \circ \text{softmax}(z),$$

denote the i -th component of $\text{softmax}(z)$, and so we calculate

$$\begin{aligned}
\frac{\partial S^i}{\partial z^j} &= \frac{\partial}{\partial z^j} \left[\left(\sum_{k=1}^m e^{z^k} \right)^{-1} e^{z^i} \right] \\
&= - \left(\sum_{k=1}^m e^{z^k} \right)^{-2} \left(\sum_{k=1}^m e^{z^k} \delta_j^k \right) e^{z^i} + \left(\sum_{k=1}^m e^{z^k} \right)^{-1} e^{z^i} \delta_j^i \\
&= - \left(\sum_{k=1}^m e^{z^k} \right)^{-2} e^{z^j} e^{z^i} + S^i \delta_j^i \\
&= -S^j S^i + S^i \delta_j^i \\
&= S^i (\delta_j^i - S^j).
\end{aligned}$$

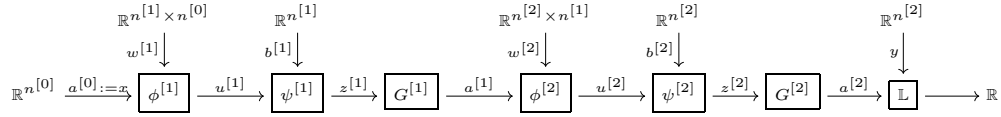
That is, as a map $dS_z : T_z \mathbb{R}^m \rightarrow T_{S(z)} \mathbb{R}^m$, we have that

$$dS_z = [S^i (\delta_j^i - S^j)]_j^i,$$

and we make note that dS_z is symmetric (i.e., it's also the reverse differential).

2.2 Backward Propagation

We consider a neural network of the form



where we have the functions:

1.

$$G^{[\ell]} : \mathbb{R}^{n^{[\ell]}} \rightarrow \mathbb{R}^{n^{[\ell]}}$$

is the broadcasting of the activation unit $g^{[\ell]} : \mathbb{R} \rightarrow \mathbb{R}$.

2.

$$\phi^{[\ell]} : \mathbb{R}^{n^{[\ell]} \times n^{[\ell-1]}} \times \mathbb{R}^{n^{[\ell-1]}} \rightarrow \mathbb{R}^{n^{[\ell]}}$$

is given by

$$\phi^{[\ell]}(w, x) = wx.$$

3.

$$\psi^{[\ell]} : \mathbb{R}^{n^{[\ell]}} \times \mathbb{R}^{n^{[\ell]}} \rightarrow \mathbb{R}^{n^{[\ell]}}$$

is given by

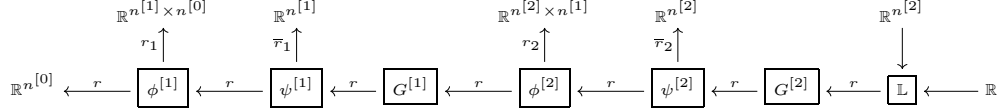
$$\psi^{[\ell]}(b, x) = x + b.$$

4.

$$\mathbb{L} : \mathbb{R}^{n^{[2]}} \times \mathbb{R}^{n^{[2]}} \rightarrow \mathbb{R}$$

is the given loss-function.

We now consider back-propagating through the neural network via “reverse exterior differentiation”. We represent our various reverse derivatives via the following diagram:



First, we need to consider our individual derivatives:

1. Suppose $G : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is the broadcasting of $g : \mathbb{R} \rightarrow \mathbb{R}$. Then for $(x, \xi) \in T\mathbb{R}^n$, we have that

$$\begin{aligned} dG_x(\xi) &= G'(x) \odot \xi \\ &= \text{diag}(G'(x)) \cdot \xi \end{aligned}$$

and for any $\zeta \in T_{G(x)}\mathbb{R}^n$, the reverse derivative is given by

$$\begin{aligned} rG_x(\zeta) &= G'(x) \odot \zeta \\ &= \text{diag}(G'(x)) \cdot \zeta. \end{aligned}$$

2. Suppose $\phi : \mathbb{R}^{m \times n} \times \mathbb{R}^n \rightarrow \mathbb{R}^m$ is given by

$$\phi(w, x) = wx.$$

Then we have two differentials to consider:

- (a) For any $(w, x) \in \mathbb{R}^{m \times n} \times \mathbb{R}^n$ and any $\xi \in T_x\mathbb{R}^n$, we have that

$$\begin{aligned} d\phi_{(w,x)}(\xi) &= w\xi \\ &= L_w(\xi); \end{aligned}$$

and for any $\zeta \in T_{\phi(w,x)}\mathbb{R}^m$, we have the reverse derivative

$$\begin{aligned} r\phi_{(w,x)}(\zeta) &= w^T \zeta \\ &= L_{w^T}(\zeta); \end{aligned}$$

where $L_A(B) = AB$, i.e., left-multiplication by A .

(b) For any $(w, x) \in \mathbb{R}^{m \times n} \times \mathbb{R}^n$ and any $\eta \in T_w\mathbb{R}^{m \times n}$ we have that

$$\begin{aligned} d_1\phi_{(w,x)}(\eta) &= \eta x \\ &= R_x(\eta); \end{aligned}$$

and for any $\zeta \in T_{\phi(w,x)}\mathbb{R}^m$, we have the reverse derivative

$$\begin{aligned} r_1\phi_{(w,x)}(\zeta) &= \zeta x^T \\ &= R_{x^T}(\zeta); \end{aligned}$$

where $R_A(B) = BA$, i.e., right-multiplication by A .

3. Suppose $\psi : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}^n$ is given by

$$\psi(b, x) = x + b.$$

Then we again have two (identical) differentials to consider:

(a) For any $(x, b) \in \mathbb{R}^n \times \mathbb{R}^n$ and any $\xi \in T_x\mathbb{R}^n$, we have that

$$d\psi_{(b,x)}(\xi) = \xi;$$

and for any $\zeta \in T_{\psi(b,x)}\mathbb{R}^n$, we have the reverse derivative

$$r\psi_{(b,x)}(\zeta) = \zeta.$$

(b) For any $(x, b) \in \mathbb{R}^n \times \mathbb{R}^n$ and any $\eta \in T_b\mathbb{R}^n$, we have that

$$d_1\psi_{(b,x)}(\eta) = \eta;$$

and for any $\zeta \in T_{(\psi(b,x))}\mathbb{R}^n$, we have the reverse derivative

$$\bar{r}_1\psi_{(b,x)}(\zeta) = \zeta.$$

Include the following two results in the Reverse Differential appendix once created.

Proposition 2.3. *Suppose we have the compositional diagram*

$$\mathbb{R}^n \xrightarrow{f} \mathbb{R}^m \xrightarrow{g} \mathbb{R}^k \xrightarrow{h} \mathbb{R}^l$$

and we let $F = h \circ g \circ f : \mathbb{R}^n \rightarrow \mathbb{R}^l$. Then for any $x \in \mathbb{R}^n$ and any $\zeta \in T_{F(x)}\mathbb{R}^l$, the reverse derivative satisfies

$$rF_x(\zeta) = rf_x \circ rg_{f(x)} \circ rh_{g(f(x))}(\zeta).$$

Proof: For any $\xi \in T_x\mathbb{R}^n$ and any $\zeta \in T_{F(x)}\mathbb{R}^l$, we have by definition

$$\begin{aligned} \langle rF_x(\zeta), \xi \rangle_{\mathbb{R}^n} &= \langle \zeta, dF_x(\xi) \rangle_{\mathbb{R}^l} \\ &= \langle \zeta, dh_{g(f(x))} \circ dg_{f(x)} \circ df_x(\xi) \rangle_{\mathbb{R}^l} \\ &= \langle rh_{g(f(x))}(\zeta), dg_{f(x)} \circ df_x(\xi) \rangle_{\mathbb{R}^k} \\ &= \langle rg_{f(x)} \circ rh_{g(f(x))}(\zeta), df_x(\xi) \rangle_{\mathbb{R}^m} \\ &= \langle rf_x \circ rg_{f(x)} \circ rh_{g(f(x))}(\zeta), \xi \rangle_{\mathbb{R}^n} \end{aligned}$$

as desired. \square

Lemma 2.4. *Suppose $f : \mathbb{R}^{n \times m} \rightarrow \mathbb{R}^k$, and for $P \in \mathbb{R}^{n \times m}$, let $R = rf_P$. Then $R \in \mathbb{R}^{k \times m}_n$ is rank $(1, 2)$ -tensor written in coordinates as*

$$R = R_i{}^\mu{}_\nu \frac{\partial}{\partial X_\nu^\mu} \otimes dx^i,$$

and the components is given by

$$R_i{}^\mu{}_\nu = \frac{\partial f^i}{\partial X_\mu^\nu}$$

Proof: Considering the basis vectors $\frac{\partial}{\partial X_\mu^\nu} \in T_P\mathbb{R}^{n \times m}$ and $\frac{\partial}{\partial x^i} \in T_{f(P)}\mathbb{R}^k$ we have that

$$\begin{aligned} R_i{}^\mu{}_\nu &= \left\langle R \left(\frac{\partial}{\partial x^i} \right), \frac{\partial}{\partial X_\mu^\nu} \right\rangle_F \\ &= \left\langle \frac{\partial}{\partial x^i}, df_P \left(\frac{\partial}{\partial X_\mu^\nu} \right) \right\rangle_{\mathbb{R}^k} \\ &= \left\langle \frac{\partial}{\partial x^i}, \frac{\partial f^\alpha}{\partial X_\mu^\nu} \frac{\partial}{\partial x^\alpha} \right\rangle_{\mathbb{R}^k} \\ &= \delta_{i\alpha} \frac{\partial f^\alpha}{\partial X_\mu^\nu}, \end{aligned}$$

as desired. \square

Returning to our neural network, for each point (x_j, y_j) in our training set, we first let

$$F_j := \mathbb{L} \circ G^{[2]} \circ \psi^{[2]} \circ \phi^{[2]} \circ G^{[1]} \circ \psi^{[1]} \circ \phi^{[1]},$$

and we have our cost function

$$\mathbb{J} := \frac{1}{N} \sum_{j=1}^N F_j.$$

We use the following notation for our inputs and outputs of our respective functions:

•

$$\phi^{[\ell]} : (w^{[\ell]}, a^{[\ell-1]}_j) \mapsto u^{[\ell]}_j,$$

•

$$\psi^{[\ell]} : (b^{[\ell]}, u^{[\ell]}_j) \mapsto z^{[\ell]}_j,$$

•

$$G^{[\ell]} : z^{[\ell]}_j \mapsto a^{[\ell]}_j.$$

Let $p = (w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]})$ is a point in our parameter space. Suppose we wish to apply gradient descent with learning rate $\alpha \in T_{\mathbb{J}(p)}\mathbb{R}$, we would define our parameter updates via

$$\begin{aligned} w^{[1]} &:= w^{[1]} - r_1 \mathbb{J}_p(\alpha) \\ b^{[1]} &:= b^{[1]} - \bar{r}_1 \mathbb{J}_p(\alpha) \\ w^{[2]} &:= w^{[2]} - r_2 \mathbb{J}_p(\alpha) \\ b^{[2]} &:= b^{[2]} - \bar{r}_2 \mathbb{J}_p(\alpha). \end{aligned}$$

Moreover, by linearity (and independence of our training data), we see that

$$r \mathbb{J}_p = \frac{1}{N} \sum_{j=1}^N r(F_j)_p,$$

so we need only calculate the various reverse derivatives of F_j .

To this end, we suppress the index j when we're working with the compositional function F . We calculate the reverse derivatives in the order traversed in our back-propagating path along the network.

1. $\bar{r}_2\mathbb{J}_p$:

$$\begin{aligned}
\bar{r}_2 F_p &= \bar{r}_2 (\mathbb{L} \circ G^{[2]} \circ \psi^{[2]})_p \\
&= \bar{r}_2 \psi_p^{[2]} \circ rG_{z^{[2]}}^{[2]} \circ r\mathbb{L}_{a^{[2]}} \\
&= \mathbb{1} \circ rG_{z^{[2]}}^{[2]} \circ r\mathbb{L}_{a^{[2]}} \\
&= rG_{z^{[2]}}^{[2]} \circ r\mathbb{L}_{a^{[2]}},
\end{aligned}$$

and hence

$$\bar{r}_2\mathbb{J}_p = \frac{1}{N} \sum_{j=1}^N rG_{z^{[2]}_j}^{[2]} \cdot r\mathbb{L}_{a^{[2]}_j}$$

2. $r_2\mathbb{J}_p$:

$$\begin{aligned}
r_2 F_p &= r_2 (\mathbb{L} \circ G^{[2]} \circ \psi^{[2]} \circ \phi^{[2]})_p \\
&= r_2 \phi_p^{[2]} \circ r\psi_{u^{[2]}}^{[2]} \circ rG_{z^{[2]}}^{[2]} \circ r\mathbb{L}_{a^{[2]}} \\
&= R_{a^{[1]}T} \circ \mathbb{1} \circ rG_{z^{[2]}}^{[2]} \circ r\mathbb{L}_{a^{[2]}} \\
&= R_{a^{[1]}T} \circ rG_{z^{[2]}}^{[2]} \circ r\mathbb{L}_{a^{[2]}},
\end{aligned}$$

and hence

$$r_2\mathbb{J}_p = \frac{1}{N} \sum_{j=1}^N R_{a^{[1]}T_j} \circ rG_{z^{[2]}_j}^{[2]} \cdot r\mathbb{L}_{a^{[2]}_j}.$$

Notice that this is not just a sum after matrix multiplication since we have composition with an operator, namely, $R_{a^{[1]}T_j}$. However, since the learning rate $\alpha \in T_{\mathbb{J}(p)}\mathbb{R} \cong \mathbb{R}$, which may pass through the aforementioned linear composition, we conclude that

$$\begin{aligned}
r_2\mathbb{J}_p &= \frac{1}{N} \sum_{j=1}^N R_{a^{[1]}T_j} \circ rG_{z^{[2]}_j}^{[2]} \cdot r\mathbb{L}_{a^{[2]}_j} \\
&= \frac{1}{N} \sum_{j=1}^N rG_{z^{[2]}_j}^{[2]} \cdot r\mathbb{L}_{a^{[2]}_j} a^{[1]T_j}.
\end{aligned}$$

3. $\bar{r}_1\mathbb{J}_p$:

$$\begin{aligned}
\bar{r}_1 F_p &= \bar{r}_1 (\mathbb{L} \circ G^{[2]} \circ \psi^{[2]} \circ \phi^{[2]} \circ G^{[1]} \circ \psi^{[1]})_p \\
&= \bar{r}_1 \psi_p^{[1]} \circ rG_{z^{[1]}}^{[1]} \circ r\phi_{a^{[1]}}^{[2]} \circ r\psi_{u^{[2]}}^{[2]} \circ rG_{z^{[2]}}^{[2]} \circ r\mathbb{L}_{a^{[2]}} \\
&= \mathbb{1} \circ rG_{z^{[1]}}^{[1]} \circ L_{w^{[2]}T} \circ \mathbb{1} \circ rG_{z^{[2]}}^{[2]} \circ r\mathbb{L}_{a^{[2]}} \\
&= rG_{z^{[1]}}^{[1]} \circ L_{w^{[2]}T} \circ rG_{z^{[2]}}^{[2]} \circ r\mathbb{L}_{a^{[2]}},
\end{aligned}$$

and hence

$$\bar{r}_1 \mathbb{J}_p = \frac{1}{N} \sum_{j=1}^N rG_{z^{[1]}_j}^{[1]} \cdot w^{[2]T} \cdot rG_{z^{[2]}_j}^{[2]} \cdot r\mathbb{L}_{a^{[2]}_j}.$$

4. $r_1 \mathbb{J}_p$:

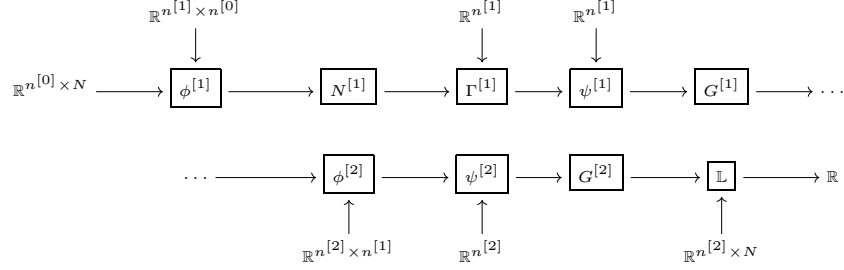
$$\begin{aligned} r_1 F_p &= r_1 (\mathbb{L} \circ G^{[2]} \circ \psi^{[2]} \circ \phi^{[2]} \circ G^{[1]} \circ \psi^{[1]} \circ \phi^{[1]})_p \\ &= r_1 \phi_p^{[1]} \circ r\psi_{u^{[1]}}^{[1]} \circ rG_{z^{[1]}}^{[1]} \circ r\phi_{a^{[1]}}^{[2]} \circ r\psi_{u^{[2]}}^{[2]} \circ rG_{z^{[2]}}^{[2]} \circ r\mathbb{L}_{a^{[2]}} \\ &= R_{x^T} \circ \mathbb{1} \circ rG_{z^{[1]}}^{[1]} \circ L_{w^{[2]T}} \circ \mathbb{1} \circ rG_{z^{[2]}}^{[2]} \circ r\mathbb{L}_{a^{[2]}} \\ &= R_{x^T} \circ rG_{z^{[1]}}^{[1]} \circ L_{w^{[2]T}} \circ rG_{z^{[2]}}^{[2]} \circ r\mathbb{L}_{a^{[2]}}, \end{aligned}$$

and hence

$$\begin{aligned} r_1 \mathbb{J}_p &= \frac{1}{N} \sum_{j=1}^N R_{x_j^T} \circ rG_{z^{[1]}_j}^{[1]} \cdot w^{[2]T} \cdot rG_{z^{[2]}_j}^{[2]} \cdot r\mathbb{L}_{a^{[2]}_j} \\ &= \frac{1}{N} \sum_{j=1}^N rG_{z^{[1]}_j}^{[1]} \cdot w^{[2]T} \cdot rG_{z^{[2]}_j}^{[2]} \cdot r\mathbb{L}_{a^{[2]}_j} \cdot x_j^T \end{aligned}$$

3 Deep Neural Networks

In this section we discuss a general “deep” neural network, which consist of L layers. That is, we have a network of the form:



3.1 Backward Propagation

As the general derivation for backpropagation can be easily (if not tediously) generalized from ?? using induction, we give the general outline for computational purposes.

Let $\mathbb{L} : \mathbb{R}^{m_L} \times \mathbb{R}^{m_L} \rightarrow \mathbb{R}$ be a generic loss function, and suppose our cost function is given by the usual

$$\mathbb{J}(W, b) = \frac{1}{n} \sum_{j=1}^n \mathbb{L}(\hat{y}_j, y_j).$$

Then from previous computations, we have the following gradients for any $\ell \in \{1, 2, \dots, L\}$, that

$$\begin{aligned} \frac{\partial \mathbb{J}}{\partial W^{[\ell]}} &= \frac{1}{n} \delta^{[\ell]} a^{[\ell-1]T} \\ \frac{\partial \mathbb{J}}{\partial b^{[\ell]}} &= \frac{1}{n} \sum_{j=1}^n \delta^{[\ell]}_j \end{aligned}$$

where we impose the notation of

$$a^{[0]} := x.$$

So we need only give a full characterization of $\delta^{[\ell]}$. To this end, we define

recursively starting at layer- L by

$$\begin{aligned}
\delta^{[L]T} &:= d(\mathbb{L}_y)_{a^{[L]}} \cdot dg_{z^{[L]}}^{[L]}, \\
\delta^{[L-1]T} &:= \delta^{[L]T} \cdot W^{[L]} \cdot dg_{z^{[L-1]}}^{[L-1]}, \\
&\vdots \\
\delta^{[\ell]T} &:= \delta^{[\ell+1]T} W^{[\ell+1]} dg_{z^{[\ell]}}^{[\ell]}, \\
&\vdots \\
\delta^{[1]T} &:= \delta^{[2]T} W^{[2]} dg_{z^{[1]}}^{[1]},
\end{aligned}$$

as desired.

3.2 Implementation in Python via numpy

We implement a neural network with an arbitrary number of layers and nodes, with the ReLU function as the activator on all hidden nodes and the sigmoid function on the output layer for binary classification with the log-loss function.

```

1  #! python3
2
3  import numpy as np
4
5  from mlLib.utils import LinearParameters, apply_activation
6
7  class NeuralNetwork():
8      def __init__(self, config):
9          """
10             Parameters:
11             -----
12             config : Dict
13                 config['lp_reg'] = 0,1,2
14                 config['nodes'] = List[int]
15                 config['bias'] = List[Boolean]
16                 config['activators'] = List[str]
17
18             Returns:
19             -----
20             None
21             """
22         self.config = config
23         self.lp_reg = config['lp_reg']
24         self.nodes = config['nodes']

```

```

25         self.bias = config['bias']
26         self.activators = config['activators']
27         self.L = len(config['nodes']) - 1
28
29     def forward_propagation(self, params, x):
30         """
31         Parameters:
32         -----
33         params : Dict[class[Parameters]]
34             params[l].w = Weights
35             params[l].bias = Boolean
36             params[l].b = Bias
37         x : array_like
38
39         Returns:
40         -----
41         cache = Dict[array_like]
42             cache['a'] = a
43             cache['dg'] = dg
44
45         """
46         # Initialize dictionaries
47         a = {}
48         dg = {}
49
50         a[0], dg[0] = apply_activation(x, self.activators[0])
51
52         for l in range(1, self.L + 1):
53             z = params[l].forward(a[l - 1])
54             a[l], dg[l] = apply_activation(z, self.activators[l])
55
56         cache = {'a' : a, 'dg' : dg}
57         return cache
58
59     def cost_function(self, params, a, y, lambda_=0.01, eps=1e-8):
60         """
61         Parameters:
62         -----
63         params: class[Parameters]
64         a: array_like
65         y: array_like
66         lambda_: float
67             Default: 0.01
68         eps: float
69             Default: 1e-8
70
71         Returns:

```

```

72         -----
73         cost: float
74         """
75         n = y.shape[1]
76         if self.lp_reg == 0:
77             lambda_ = 0.0
78
79         # Compute regularization term
80         R = 0
81         for param in params.values():
82             R += np.sum(np.abs(param.w) ** self.lp_reg)
83         R *= (lambda_ / (2 * n))
84
85         # Compute unregularized cost
86         a = np.clip(a, eps, 1 - eps) # Bound a for stability
87         J = (-1 / n) * (np.sum(y * np.log(a) + (1 - y) * np.log(1 - a)))
88
89         cost = float(np.squeeze(J + R))
90
91         return cost
92
93     def backward_propagation(self, params, cache, y):
94         """
95         Parameters:
96         -----
97         params : Dict[class[Parameters]]
98                 params[l].w = Weights
99                 params[l].bias = Boolean
100                params[l].b = Bias
101         cache : Dict[array_like]
102                 cache['a'] : array_like
103                 cache['dg'] : array_like
104         y : array_like
105
106         Returns:
107         -----
108         None
109         """
110
111         # Retrieve cache
112         a = cache['a']
113         dg = cache['dg']
114
115         # Initialize differentials along the network
116         delta = {}
117         delta[self.L] = (a[self.L] - y) / y.shape[1]
118

```

```

119         for l in reversed(range(1, self.L + 1)):
120             delta[l - 1] = dg[l - 1] * params[l].backward(delta[l], a[l - 1])
121
122     def update_parameters(self, params, learning_rate=0.1):
123         """
124         Parameters:
125         -----
126         params : Dict[class[Parameters]]
127             params[l].w = Weights
128             params[l].bias = Boolean
129             params[l].b = Bias
130         learning_rate : float
131             Default : 0.01
132
133         Returns:
134         -----
135         None
136         """
137         for param in params.values():
138             param.update(learning_rate)
139
140     def fit(self, x, y, learning_rate=0.1, lambda_=0.01, num_iters=10000):
141         """
142         Parameters:
143         -----
144         x : array_like
145         y : array_like
146         learning_rate : float
147             Default : 0.1
148         lambda_ : float
149             Default : 0.0
150         num_iters : int
151             Default : 10000
152
153         Returns:
154         -----
155         costs : List[floats]
156         params : class[Parameters]
157         """
158         # Initialize parameters per layer
159         params = {}
160         for l in range(1, self.L + 1):
161             params[l] = LinearParameters((self.nodes[l], self.nodes[l - 1]), self.b)
162
163         costs = []
164         for i in range(num_iters):
165             cache = self.forward_propagation(params, x)

```

```

166         cost = self.cost_function(params, cache['a'][self.L], y, lambda_)
167         costs.append(cost)
168         self.backward_propagation(params, cache, y)
169         self.update_parameters(params, learning_rate)
170
171         if i % 1000 == 0:
172             print(f'Cost_after_iteration_{i}:_{cost}')
173
174     return params
175
176 def evaluate(self, params, x):
177     """
178     Parameters:
179     -----
180     params : class[Parameters]
181     x : array_like
182
183     Returns:
184     -----
185     y_hat : array_like
186     """
187     cache = self.forward_propagation(params, x)
188     a = cache['a'][self.L]
189     y_hat = (~(a < 0.5)).astype(int)
190     return y_hat
191
192 def accuracy(self, params, x, y):
193     """
194     Parameters:
195     -----
196     params : class[Parameters]
197     x : array_like
198     y : array_like
199
200     Returns:
201     -----
202     accuracy : float
203     """
204     y_hat = self.evaluate(params, x)
205     acc = np.sum(y_hat == y) / y.shape[1]
206
207     return acc

```

3.3 Implementation in Python via tensorflow

We implement a neural network using tensorflow.keras.

```

1 #! python3
2
3 import pandas as pd
4 import numpy as np
5 from sklearn.model_selection import train_test_split
6 from tensorflow import keras
7 from keras import Model, Input
8 from keras.layers import Dense
9
10 def keras_functional_nn(csv):
11     df = pd.read_csv(csv)
12     dataset = df.values
13     x, y = dataset[:, :-1], dataset[:, -1].reshape(-1, 1)
14     x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.15)
15     train = {'x' : x_train, 'y' : y_train}
16     test = {'x' : x_test, 'y' : y_test}
17     mu = np.mean(train['x'], axis=0, keepdims=True)
18     var = np.var(train['x'], axis=0, keepdims=True)
19     train['x'] = (train['x'] - mu) / np.sqrt(var)
20     test['x'] = (test['x'] - mu) / np.sqrt(var)
21
22     ## Define network structure
23     input_layer = Input(shape=(10,))
24     hidden_layer_1 = Dense(
25         32,
26         activation='relu',
27         kernel_initializer='he_normal',
28         bias_initializer='zeros'
29     )(input_layer)
30     hidden_layer_2 = Dense(
31         8,
32         activation='relu',
33         kernel_initializer='he_normal',
34         bias_initializer='zeros'
35     )(hidden_layer_1)
36     output_layer = Dense(
37         1,
38         activation='sigmoid',
39         kernel_initializer='he_normal',
40         bias_initializer='zeros'
41     )(hidden_layer_2)
42
43     model = Model(inputs=input_layer, outputs=output_layer)
44     model.summary()
45
46     ## Compile desired model
47     model.compile(

```

```

48         loss='binary_crossentropy',
49         optimizer='adam',
50         metrics=['accuracy']
51     )
52
53     ## Train the model
54     hist = model.fit(
55         train['x'],
56         train['y'],
57         batch_size=32,
58         epochs=150,
59         validation_split=0.17
60     )
61
62     ## Evaluate the model
63     test_scores = model.evaluate(test['x'], test['y'], verbose=2)
64     print(f'Test Loss: {test_scores[0]}')
65     print(f'Test Accuracy: {test_scores[1]}')

```


References