# Deep Neural Networks - A bird's-eye view

**Jose Luis Silva**[*]
Department of Physics and Astronomy
Uppsala University
Uppsala, SE 751 20
`jose.silva@physics.uu.se`

## Abstract

********* to be updated ***********
In this report we will explore the logistic regression method for binary classification. As part of the assignments, I will also provide details about the implementations in Python. This means an explanatory bird's-eye view of how to optimize parameters of the model using gradient descent method, Binary Cross-Entropy loss function and cost function derivatives. Additionally, we will use the biopsy breast cancer dataset to deploy the model and evaluate the accuracy to predict if a certain biopsy is "benign" or "malignant".

## 1 Introduction

********* to be updated ***********
In binary classification tasks, our main interest is to determine optimized parameters of a model that can efficiently predict the probability mass function $p(y = 1|x)$ of an event $y \in \{0, 1\}$ given x on a specific interval $[0, 1]$. In our case, we want to find a set parameters using the logistic regression model to predict the probability vector $p_i = P(y_i = 1|x_i)$ for a data set with $i = \{1, 2, ..., n - 1, n\}$ inputs vectors $\mathbf{x}_i$ and $y_i$ outputs where $y_i \in \{0, 1\}$, $\mathbf{x}_i = [x_{i1}, x_{i2}, .., x_{i(p-1)}, x_{ip}]^T$ and $p = \{1, 2, ..., m - 1, m\}$ features. For each input vector $\mathbf{x}_i$, logistic regression model can be summarized as computing the following set of linear combinations:

$$
\begin{aligned}
z_1 &= \sum_{j=1}^{p} w_j x_1 j + b \\
z_2 &= \sum_{j=1}^{p} w_j x_2 j + b \\
&\vdots \\
z_n &= \sum_{j=1}^{p} w_j x_n j + b = \mathbf{w}^T \mathbf{x}_i + b
\end{aligned}
\tag{1}
$$

where $n$ is the number of input vectors, $p$ the number of features in each input vector, $b$ represents the offset parameter and $\mathbf{w} = [w_1, w_2, ..., w_p]^T$ the weights. As a first step we can generate a set of normalized random values to initialize the weights and the offset b. This procedure allows the computation of the probability that each input vector $\mathbf{x}_i$ belongs to certain class $y_i = \{1, 0\}$. Hence,

---

[*]https://jseluis.github.io/

we can apply the activation function(sigmoid) over $\mathbf{z}_i$, such that:

$$p_i = \sigma(z_i) = \left[1 + \exp\left(-\sum_{j=1}^{p} w_j x_{ij} + b\right)\right]^{-1}$$

(2)

This procedure is followed by an optimization of the model through the maximization of the likelihood for a given set of parameters. We use the cross-entropy loss $L(p_i, y_i)$ function for each input vector and take an average over the whole data set to estimate an initial cost function $J(p_i, y_i)$, where:

$$J = \frac{1}{n}\sum_{i=1}^{n} L(p_i, y_i)$$

(3)

$$L(p_i, y_i) = -y_i \ln(p_i) - (1 - y_i)\ln(1 - p_i)$$

(4)

However, the parameters of the model are optimized through a stochastic gradient descent method, which makes use of the gradient of the cost function with respect to the parameters of the model. We need to update the weights $w_i$ and the offset b for , such that:

$$w_1^l = w_1^{(l-1)} - \alpha\frac{\partial J}{\partial w_1}$$

$$w_2^l = w_2^{(l-1)} - \alpha\frac{\partial J}{\partial w_2}$$

$$\vdots$$

$$w_n^l = w_n^{(l-1)} - \alpha\frac{\partial J}{\partial w_n}$$

(5)

$$b^l = b^{(l-1)} - \alpha\frac{\partial J}{\partial b}$$

(6)

where (l-1) represents parameter from the previous iteration and $\alpha$ is responsible to modulate the learning rate. We need to determine the expressions for $\frac{\partial J}{\partial w_n}$ and $\frac{\partial J}{\partial b}$ as a function of the input data and parameter of the model. By using the chain rule, we can derive both equations as follows:

$$\frac{\partial J}{\partial w_j} = \frac{\partial J}{\partial z_i}\frac{\partial z_i}{\partial w_j} = \frac{\partial J}{\partial z_i}\left(\frac{\partial}{\partial w_j}\right)\left(\sum_{i=1}^{p} w_j x_{ij} + b\right) = \left(\frac{\partial J}{\partial z_i}\right)x_{ij}$$

(7)

$$\frac{\partial J}{\partial b} = \frac{\partial J}{\partial z_i}\frac{\partial z_i}{\partial b} = \frac{\partial J}{\partial z_i}\left(\frac{\partial}{\partial b}\right)\left(\sum_{i=1}^{p} w_j x_{ij} + b\right) = \left(\frac{\partial J}{\partial z_i}\right)$$

(8)

From equation (7) and (8), $\frac{\partial J}{\partial z_i}$ can be estimated using the derivation chain rule with the partial derivative of the cross-entropy loss function in respect to $p_i$, as follows:

$$\frac{\partial J}{\partial z_i} = \frac{1}{n}\sum_{i=1}^{n}\frac{\partial L_i}{\partial z_i} = \frac{1}{n}\sum_{i=1}^{n}\frac{\partial L_i}{\partial p_i}\frac{\partial p_i}{\partial z_i}$$

(9)

such that:

$$\frac{\partial p_i}{\partial z_i} = \left(\frac{\partial}{\partial z_i}\right)\left[1 + \exp\left(-\sum_{j=1}^{p} w_j x_{ij} + b\right)\right]^{-1} = e^{-z_i}(1 + e^{-z_i})^{-2}$$

(10)

and since:

$$p_i = (1 + e^{-z_i})^{-1} \rightarrow e^{-z_i} = p^{-1} - 1 = \frac{(1 - p_i)}{p_i}$$

(11)

equation (10) becomes:

$$\frac{\partial p_i}{\partial z_i} = \frac{(1 - p_i)(1 + e^{-z_i})^{-2}}{p_i} = p_i(1 - p_i)$$

(12)

2

The first partial derivative from equation (9) becomes:

$$\frac{\partial L_i}{\partial p_i} = -\frac{y_i}{p_i} + \frac{(1 - y_i)}{(1 - p_i)} \tag{13}$$

therefore:

$$\frac{\partial J}{\partial z_i} = \frac{1}{n} \sum_{i=1}^{n} \left[ -\frac{y_i}{p_i} + \frac{(1 - y_i)}{(1 - p_i)} \right] \left[ \frac{(1 - p_i)}{p_i} \right] = \frac{1}{n} \sum_{i=1}^{n} [-y_i(1 - p_i) + p_i(1 - y_i)]$$

$$= \frac{1}{n} \sum_{i=1}^{n} (p_i - y_i) \tag{14}$$

In order to update the parameters of the model, the following set of recurrent formulas were implemented in python:

$$w_1^l = w_1^{(l-1)} - \frac{\alpha}{n} \sum_{i=1}^{n} (p_i - y_i) x_{1j}$$

$$w_2^l = w_2^{(l-1)} - \alpha \frac{\alpha}{n} \sum_{i=1}^{n} (p_i - y_i) x_{2j}$$

$$\vdots$$

$$w_n^l = w_n^{(l-1)} - \frac{\alpha}{n} \sum_{i=1}^{n} (p_i - y_i) x_{nj} \tag{15}$$

$$b^l = b^{(l-1)} - \frac{\alpha}{n} \sum_{i=1}^{n} (p_i - y_i) \tag{16}$$

For the exercise 1.3, the logistic regression model that can be trained with stochastic gradient descent was implemented in Python using Numpy library for vectorization and faster operations. I have created following trained data set: $n = 6$ input vectors with two features $x_i = [x_1, x_2]$ and outputs $y_i = 1, 0$ where $y_i = 1$ if $x_1 = x_2$ or $x_1 > x_2$ and $y_i = 0$ if $x_2 > x_1$:

$$x\_train = np.array([[1, 1], [3, 4], [5, 5], [7, 7], [1, 4], [4, 4]]) \tag{17}$$
$$x\_test = np.array([[3, 3], [3, 10], [1, 5], [12, 12], [1, 0], [10, 100]]) \tag{18}$$
$$y\_train = np.array([1, 0, 1, 1, 0, 1]) \tag{19}$$
$$y\_test = np.array([1, 0, 0, 1, 1, 0]) \tag{20}$$
$$\tag{21}$$

An initialization function was responsible for starting a set of parameter such as a fixed learning rate of 0.5 and break of the self-consistent loops based on values of the cost function which should be lower than 0.001 (convergence break). The training was performed for 18151 steps until the break of the while loop due to convergence. The set of equations for the sigmoid activation function, cost function and proper deduced gradients to update the parameters of the model during every iteration were implemented as functions in Python and called inside the while loop. The accuracy for the predictions based on the optimized parameters was 100% and consequently the confusion matrix matched the number of True Positives and True Negatives of the test data set. The code is well explained in the appendix session.

For the exercise 1.4, I have used the same python functions and part of the previous implementation with some changes in order to predict the biopsy as malignant or benign in the test data set using the optimized model. First we initialized the variables and used a learning rate of 0,8. We reshaped Y for train and test to (300,) and (383,) respectively. Additionally I have applied a normalization of $X_{train}$ and $X_{test}$ dividing by the maximum value of the features. The functions to estimate the probabilities, cost function and gradients are optimized for any $X_{train}$ and $X_{train}$ shape. I have included arrays to

3

append the cost function during the optimization of the parameters for both training and test data set in order to plot the cost function x number of iterations. We optimized the parameters with 4577 iterations considering a threshold $min_{cost}$ of 0.61. The cost function decays very fast as a function of the iteration, as shown in the Figure 1.
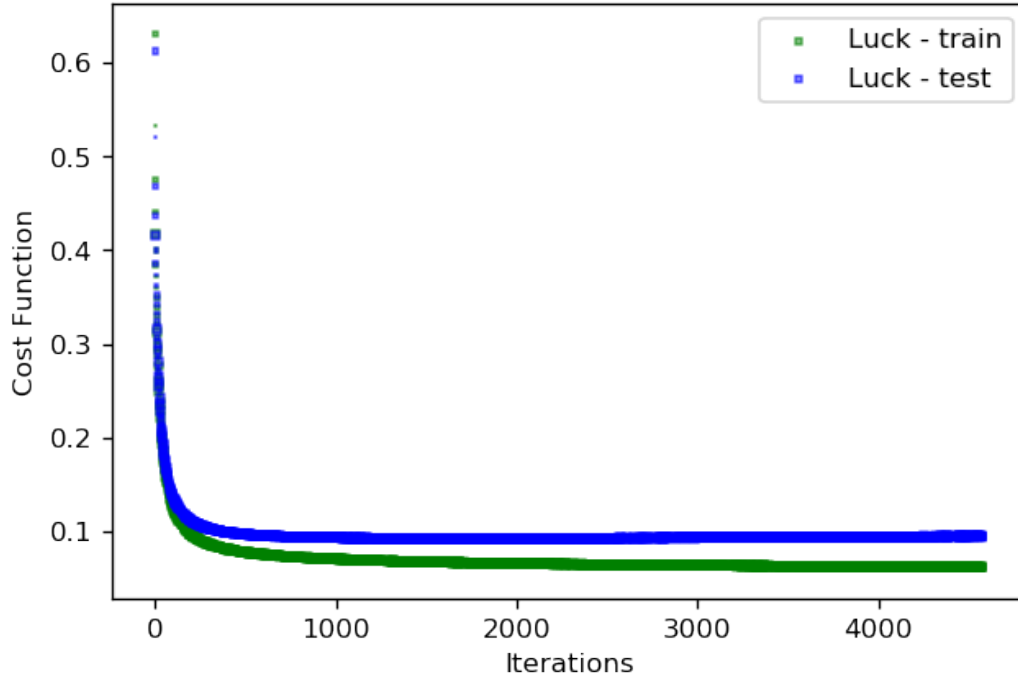


Figure 1: Cost function x Iteration - Blue = Test data set , Green = Train data set

As expected, the the cost function is lower for the training data set (green color) as compared to the testing data set. The accuracy of the model on the training and test data set was 0.97 and 0.96, respectively. The confusion matrix for the training and test data set showed a very low number of False Negatives with 5 missed predictions as malignant in both cases. We conclude that the logistic regression model optimized with stochastic gradient descent is very accurate for binary classification predictions in this particular data set.

## 2 Appendix

Exercise 1.3

```python
import numpy as np
import matplotlib.pyplot as plt

#   The shape of x_train and x_test is (6, 2) and  y_test and y_train is
#    (6,). Each vector has two features x1,x2 that belongs to a class y=0
#    or y=1. If x1=x2 or #  x1>x2 then y = 0 otherwise y=0.

x_train = np.array([[1,1],[3,4],[5,5],[7,7],[1,4],[4,4]])
y_train = np.array([1,0,1,1,0,1])
x_test = np.array([[3,3],[3,10],[1,5],[12,12],[1,0],[10,100]])
y_test = np.array([1,0,0,1,1,0])

#   Initialization of the parameters through initialize() function. The
#    weights were chosen randomly using a normal distribution with
#    gaussians. The transpose of x_train is initialized for future matrix
#    operations. The initial offset b = 0 and min_cost controls the number
#     of iterations to minimize the cost function by breaking the self-
#    consistent loop when the cost function is lower than 0.001

def initialize(x_train, y_train):
    w = np.random.random([x_train.shape[1]])
    n = x_train.shape[0]
    x = np.transpose(x_train)
    b = 0
    alpha = -0.5
    min_cost = 0.001 # Control variable to stop the minimization of the
    cost function
    return w,x,n,b,alpha,min_cost

#   Definition of activation function sigma(z) based on the sigmoid. the
#    function needs to be called by defining the type of activation. i.e
#    sigma(z,activation='sigmoid'). This is important if we explore
#    possible different activation functions in the next assignments.

def sigma(z_i, activation=False):
    if(activation==False):
        return print('Please choose an activation function')
    elif(activation=='sigmoid'):
        sig = 1/(1+np.exp(-z_i))
        return sig

#   This function calculates the cross-entropy loss function and the cost
#     function. Instead of using for-loops element-wise operations are
#    performed with arrays. You need to define which type of loss function
#     to use: cost_function(y_train,p_i,loss_function='cross_entropy')

def cost_function(y_i, p_i, loss_function=False):
    if(loss_function==False):
        return print('Please choose a loss function')
    elif(loss_function=='cross_entropy'):
        n = y_i.shape[0]
        loss_calc = -(y_i*np.log(p_i)+ (1-y_i)*(np.log(1-p_i)))
        cost_calc = (1/n)*np.sum(loss_calc)
        return loss_calc,cost_calc

#   This function calculates the gradients and performs matrix operations
#     to estimate the partial derivatives. dLdb represents the derivative
#    of the loss function with respect to the offset parameter b. dJdb =
#    derivative of the cost function with respect to b and dJdw =
#    derivative of the cost function with respect to a specific j-th
```

weight. Numpy library was used for vectorization and to perform element−wise operations instead of for loops.

```python
43
44  def grad(p_in,y_in,x_in): # input p_in = array of probability, y_in =
         array of classes, x_in = matrix array.
45      n_in = y_in.shape[0]
46      dLdb = p_in − y_in
47      dJdb = (1/n_in)*np.sum(dLdb) # Derivative of the Cost function with
         respect to offset
48      dJdw = (1/n)*np.dot(dLdb,x_in) # Vectorization to facilitate the
         operations and automatically perform the sum.
49      return dJdw,dJdb
50
51  #    Initialize w=weight, x=input vectors, n= number of input vectors in
         the data set, alpha = learning rate, min_cost = control variable
52
53  w,x,n,b,alpha,min_cost=initialize(x_train,y_train)
54
55  #    Calculate an initial z using matrix formulation where  x =
         transpose of x_train. Here we also calculate the sigmoid for each z
         and print the values.
56
57  z = np.dot(w,x)+b
58  p_i=sigma(z,activation='sigmoid')
59  print('z =',z,'\ny_i =',y_train,'n =',n, 'w =',w)
60  print('p_i = ',p_i)
61
62  #   Function call to estimate the cost function based on the cross−
         entropy loss function and print both the cost function and the loss
         function. l = loss function, j = cost function
63
64  l,j= cost_function(y_train,p_i,loss_function='cross_entropy')
65  print('Cost Function =',j,'\nLoss Function vector =',l)
66
67  #   Function call to estimate the gradients dw = dJdw and db = dJdb.
         These values are used to update the parameter of the model.
68
69  dw,db= grad(p_i,y_train,x_train)
70  print(dw,db)
71
72  #   First update of the weight and offset parameters.
73
74  w += alpha*dw
75  b += alpha*db
76  print('w =',w,'\n b =',b)
77
78  #   Definition of a counter and threshold of min_cost = 0.001 for
         convergence. This means that the parameters of the model are
         optimized. All the previous operations are performed inside the loop
         until the break of the self−consistent while−loop.
79
80  counter = 0
81  while(j>min_cost):
82      z = np.dot(w,x)+b
83      counter+=1
84      p_i=sigma(z,activation='sigmoid')
85      l,j= cost_function(y_train,p_i,loss_function='cross_entropy')
86      print('Cost Function =',j,'\n')
87      dw,db = grad(p_i,y_train,x_train)
88      w +=alpha*dw
89      b +=alpha*db
90  print('Number of iterations = ',counter)
91
```

```
92 #    In order to test the optimization over the training data set,
          prediction_train checks each probability found with the new
          parameters and approximates to 1 if p_i > 0.5 and 0 otherwise.
93
94 prediction_train = np.where(p_i >0.5,1,0)
95 print(prediction_train)
96
97 #    Calculate the prediction over the test dataset
98 z_test = np.dot(w,np.transpose(x_test))+b
99 p_i_test=sigma(z_test,activation='sigmoid') # probability using optimized
          parameters
100 prediction_test = np.where(p_i_test >0.5,1,0)   # threshold of 0.5 to
          approximate y values
101 print(prediction_test)
102 np.mean(prediction_test == y_test) # Calculate the accuracy
103 print(pd.crosstab(prediction_test, y_test)) # Check the confusion matrix
```

### Exercise 1.4

```
1 # Import Numpy, Matplotlib and Pandas library
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import pandas as pd
5
6 #    Function to initialize the biopsy data set.
7
8 def load_biopsy():
9       # import data
10      biopsy = pd.read_csv('biopsy.csv', na_values='?',
11                           dtype={'ID': str}).dropna().reset_index()
12
13      # Split in training and test data
14      trainI = np.random.choice(biopsy.shape[0], size=300, replace=False)
15      trainIndex=biopsy.index.isin(trainI)
16      train=biopsy.iloc[trainIndex] # training set
17      test=biopsy.iloc[~trainIndex] # test set
18
19      # Extract relevant data features
20      X_train = train[['V1','V2','V3','V4','V5','V6','V7','V8','V9']].
          values
21      X_test = test[['V1','V2','V3','V4','V5','V6','V7','V8','V9']].values
22      Y_train=(train['class']=='malignant').astype(int).values.reshape
          ((-1,1))
23      Y_test=(test['class']=='malignant').astype(int).values.reshape((-1,1)
          )
24
25      return X_train, Y_train, X_test, Y_test
26
27 # Initialization of the parameters through initialize() function.
28 # The weights were chosen randomly using a normal distribution with
          gaussians.
29 # The transpose of x_train is initialized for future matrix operations.
30 # The initial offset b = 0 and min_cost  controls the number of
          iterations to minimize the cost function by breaking the self-
          consistent loop when the cost function is lower than 0.062
31
32 def initialize(x_train,y_train):
33      w = np.random.random([x_train.shape[1]])
34    # w=np.array([0.39401661, 0.47478989, 0.06309985, 0.99740423,
          0.33530285,0.60437357, 0.74371789, 0.3407668 , 0.81388953]) # Initial
           parameters generated randomly.
35      n = x_train.shape[0]
36      x = np.transpose(x_train)
37      b = 0
38      alpha = -0.8 # Learning rate of 0.8
```

```python
39        min_cost = 0.062
40        return w,x,n,b,alpha,min_cost
41
42  # Definition of activation function sigma(z) based on the sigmoid.
43  # This function needs to be called by defining the type of activation
         function. i.e sigma(z,activation='sigmoid').
44  # This is important if we explore possible different activation functions
          in the next assignments.
45
46  def sigma(z_i,activation=False):
47      if(activation==False):
48          return print('Please choose an activation function')
49      elif(activation=='sigmoid'):
50          sig = 1/(1+np.exp(-z_i))
51          return sig
52
53  # This function calculates the cross-entropy loss function and the cost
         function.
54  # Instead of using for-loops element-wise operations are performed with
         arrays.
55  # You need to define which type of loss function to use: cost_function(
         y_train,p_i,loss_function='cross_entropy')
56
57  def cost_function(y_i,p_i,loss_function=False):
58      if(loss_function==False):
59          return print('Please choose a loss function')
60      elif(loss_function=='cross_entropy'):
61          n = y_i.shape[0]
62          loss_calc = -(y_i*np.log(p_i) + (1-y_i)*(np.log(1-p_i)))
63          cost_calc = (1/n)*np.sum(loss_calc)
64          return loss_calc,cost_calc
65
66  # This function calculates the gradients and performs matrix operations
         to estimate the partial derivatives.
67  # dLdb represents the derivative of the loss function with respect to the
          offset parameter b.
68  # dJdb = derivative of the cost function with respect to b and dJdw =
         derivative of the cost
69  # function with respect to a specific j-th weight.
70  # Numpy library was used for vectorization and to perform element-wise
         operations instead of for loops.
71
72  def grad(p_in,y_in,x_in):
73      n_in = y_in.shape[0]
74      dLdb = p_in - y_in
75      dJdb = (1/n_in)*np.sum(dLdb)
76      dJdw = (1/n)*np.dot(dLdb,x_in) # vector dJ/dW_j to update w_j
77      return dJdw,dJdb
78
79  # Load biopsy data in training and testing variables
80  X_train, Y_train, X_test, Y_test=load_biopsy()
81  # Normalization of the training data set
82  X_train = X_train/np.max(X_train)
83  Y_train = np.reshape(Y_train,(Y_train.shape[0],)) # Reshape of the
         Y_train from (300, 1) to (300,)
84   # Normalization of the testing data set
85  X_test = X_test/np.max(X_test)
86  Y_test = np.reshape(Y_test,(Y_test.shape[0],)) # Reshape of the Y_train
         from (300, 1) to (300,)
87
88  # Initialize parameters and transpose of X_train.
89  w,x,n,b,alpha,min_cost=initialize(X_train,Y_train)
90
91  # Initialize vector z and the probabilities p_i with the initial
         parameter of the model
```

```python
92  z = np.dot(w,x)+b
93  p_i=sigma(z,activation='sigmoid')
94
95  # Function call to estimate the cost function based on the cross-entropy
        loss function and print both the cost function and the loss function.
        l = loss function, j = cost function
96  l,j= cost_function(Y_train,p_i,loss_function='cross_entropy') # l = loss
        function
97  print('Cost Function =',j,'\nLoss Function vector =',l)
98
99  # Function call to estimate the gradients dw = dJdw and db = dJdb.
100 # These values are used to update the parameter of the model.
101 dw,db= grad(p_i,Y_train,X_train)
102 print(dw,db)
103 w += alpha*dw # update weights  gradient descent
104 b += alpha*db # update the offset using gradient descent
105 print('w =',w,'\n b =',b) # print new parameters
106
107 # Definition of a counter and threshold of min_cost for convergence. This
        means that the parameters of the model are optimized. All the
        previous operations are performed inside the loop until the break of
        the self-consistent while.
108
109 counter = 0 # counter the iterations
110 cost = [] # append cost function of the training data set in the array
111 iter = [];  # append iteration in the array
112 cost_test = []; # append cost function of the test data set in the array
113 while(j>min_cost):
114     z = np.dot(w,x)+b
115     counter+=1
116     p_i=sigma(z,activation='sigmoid')
117     l,j= cost_function(Y_train,p_i,loss_function='cross_entropy')
118     cost.append(j)
119     iter.append(counter)
120
121     # calculate z, p_i and cost function based on optimized parameters
        for the test data set.
122     z_test = np.dot(w,np.transpose(X_test))+b
123     p_i_test=sigma(z_test,activation='sigmoid')
124     l_test,j_test= cost_function(Y_test,p_i_test,loss_function='
        cross_entropy')
125     cost_test.append(j_test)
126
127      print('Cost Function =',j,'\n')
128     dw,db = grad(p_i,Y_train,X_train)
129     w +=alpha*dw
130     b +=alpha*db
131 print('Number of iterations = ',counter)
132
133 prediction_train = np.where(p_i>0.5,1,0)
134
135 np.mean(prediction_train == Y_train) # Accuracy of 0.98 on training data
        set
136
137 print(pd.crosstab(prediction_train, Y_train))  # Generate the confusion
        matrix
138
139 # generate z and probabilities based on the optimized parameters
140 z_test = np.dot(w,np.transpose(X_test))+b
141 p_i_test=sigma(z_test,activation='sigmoid')
142 prediction_test = np.where(p_i_test>0.5,1,0)
143
144 # Predict the acurracy -  0.9660574412532638 on the test data set
145 np.mean(prediction_test == Y_test)
146
```

```
147 print(pd.crosstab(prediction_test, Y_test)) # Generate the confussion
         matrix for further analysis
148
149 # Plot the Cost Function x Iteration during the optimization for the
         training and test data set
150
151 plt.figure(dpi = 120)
152 s = np.random.rand(*x.shape)*10
153 plt.scatter(iter,cost,s,c="g", marker="s", alpha=0.5,label="Luck − train"
         )
154 plt.scatter(iter,cost_test,s,c="b", marker="s", alpha=0.5,label="Luck −
         test")
155 plt.xlabel("Iterations")
156 plt.ylabel("Cost Function")
157 plt.legend(loc='upper right')
158 plt.show()
```

## 3   References

******** to be updated ***********