# Classification of two class problem using Logistic Regression in python from scratch

**Sankeerth Tella**
50317364
*stella3@buffalo.edu*

## Abstract

In this project, we perform classification of two class problem using logistic regression technique. The dataset provided is Wisconsin diagnostic breast cancer which is basically classified into two classes namely Benign and Malignant. So by using logistic regression technique I developed a model of two classes so that for given input it maps into either of the classes.

## 1    INTRODUCTION

In the given task we have cancer dataset and it is basically a two class problem which comprises of two classes namely benign and malignant. We need to classify and develop a model using logistic regression. Before coming to logistic regression let's see under what category of learning strategy this logistic regression falls under. Based on the way how they learn or train a model they are classified into four categories

    i)      supervised learning
    ii)     unsupervised learning
    iii)    semi-supervised learning
    iv)    reinforcement learning

In our given scenario logistic regression falls under supervised learning. Supervised learning is nothing but model is trained on a labelled dataset. Labelled dataset is the dataset which have both input and output parameters.

Logistic regression is the classification algorithm used to develop a model and assign the observations to any of the discrete classes. Unlike linear regression which is used when dependent variable is continuous, logistic regression produces output using sigmoid function to return a probability value and that probability value is mapped to any one of the discrete classes.

### 1.1    LOGISTIC REGRESSION

Logistic regression is one of the popular supervised machine learning algorithm especially for binary classification. In simple words logistic regression is nothing but estimating parameters using logistic model and mapping them to various classes. Going in deep based on the output classifications logistic regression is further divided into three types

    i)      Binary logistic regression:

Binary logistic regression is nothing but it has only two possible outcomes and there will be only two classes to classify. Consider simple example of student data set where we have the student marks and their respective grades. Based on the given data we need to classify it into two classes namely pass and fail. So for any given new input our designed model must map to appropriate class.

    ii)     Multinomial logistic regression

In this type of regression we generalize our model to multiclass problem i.e with more than two possible outcomes. For example consider animal data set where we have characteristics

47    of more then 2 animals and we need to develop model using this given features. For any given
48    input we need to map to appropriate class.

49       iii)      Ordinal logistic regression

50    Ordinal logistic regression is a statistical technique that is used to predict behavior of ordinal
51    level dependent variables with a set of independent variables.

52

53    **1.2**      **CORE OF LOGISTIC REGRESSION**

54    **1.2.1**      **SIGMOID FUNCTION**

55    Logistic regression is named for the function used at the core of the method, the logistic
56    function. In order to map predicted values to probabilities, we use this sigmoid function. This
57    function simply maps any given value to scale between 0 and 1.

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

58
59
60    **2**      **DATASET**

61
62    Wisconsin Diagnostic Breast Cancer (WDBC) dataset is used for training, validation and testing.
63    The dataset contains 569 instances with 32 attributes (ID, diagnosis (B/M), 30 real-valued input
64    features). Features are computed from a digitized image of a fine needle aspirate (FNA) of a breast
65    mass. Computed features describes the following characteristics of the nuclei present in the image.

66

| 1 | radius (mean of distances from center to points on the perimeter) |
|---|---|
| 2 | texture (standard deviation of gray-scale values) |
| 3 | perimeter |
| 4 | area |
| 5 | smoothness (local variation in radius lengths) |
| 6 | compactness ($perimeter^2/area - 1.0$) |
| 7 | concavity (severity of concave portions of the contour) |
| 8 | concave points (number of concave portions of the contour) |
| 9 | symmetry |
| 10 | fractal dimension ("coastline approximation" - 1) |

67
68
69    The mean, standard error, and \worst" or largest (mean of the three largest values) of these features
70    were computed for each image, resulting in 30 features.

71
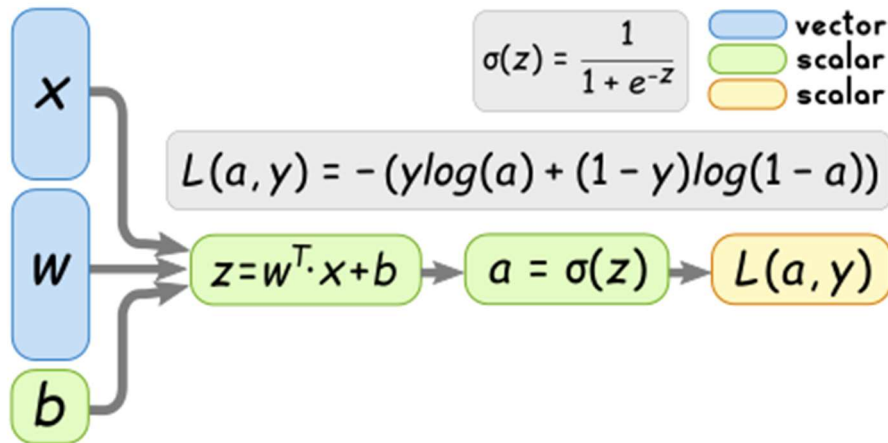72    **3**      **PREPROCESSING**

73    Before processing the given data there needs to be some preprocessing done to the data.
74    Initially I loaded the entire data using readcsv function imported from pandas library. Then
75    from whole data I selected $2^{nd}$ column which is of output column and stored in y. Next I stored
76    all the remaining columns other than y and ids i.e first column into x. Later I splitted the whole
77    data into three parts namely training data, testing data, validation data. Training data comprises
78    of 80% of given data and remaining data is divided between testing and validation 10% each.
79    And next after splitting the data scaling of data is done using standardscalar function which is
80    imported from sklearn kit.
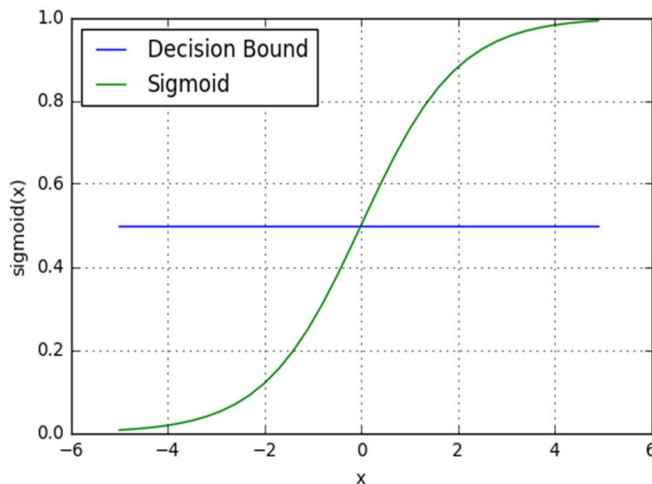
# 4 ARCHITECTURE

## 4.1 COMPUTATIONAL GRAPH

The computational graph of the logistic regression can be visualized as follows:



w, x are input vectors and their size depends on input variables.

## 4.2 DECISION BOUNDARY

After calculating the probabilistic value we need to map it to either of the class. So we need set threshold and based upon the threshold value we will divide the classes.



## 4.3 COST FUNCTION

Instead of Mean Squared Error, we use a cost function called cross-entropy, also known as Log Loss. Cross-entropy loss can be divided into two separate cost functions: one for $y=1$ and

101     one for y=0.

102

$$J(\theta) = \frac{1}{m} \sum_{i=1}^{m} \text{Cost}(h_\theta(x^{(i)}), y^{(i)})$$

$$\text{Cost}(h_\theta(x), y) = -\log(h_\theta(x)) \qquad \text{if } y = 1$$
$$\text{Cost}(h_\theta(x), y) = -\log(1 - h_\theta(x)) \qquad \text{if } y = 0$$

103

104

## 4.4     GRADIENT DESCENT

106 The objective of gradient descent is to find out optimal parameters that result in optimising a
107 given cost function. In the Logistic Regression algorithm, the optimal parameters are found
108 by minimizing the loss function and updating the weights and bias.

109



110

111

112

$$\frac{\partial J}{\partial w} = \frac{1}{m}\left\{-b\left(1-a\right) + \left(1-b\right)a\right) x\right\}$$

$$\frac{\partial J}{\partial w} = \frac{1}{m}\left\{(a-b)\; x\right\}$$

weights are new

$$w_i = w_i - \alpha\frac{\partial L}{\partial w_i}$$

Similarly same process is applied for bias

$$b = b - \alpha\frac{\partial L}{\partial b}$$

113
114

115 After minimizing the cost function through various epochs, we get the updated weights and
116 bias. And using these weights and bias we predict the output values again using sigmoid
117 function with updated final weights and bias. After predicting the final values, we need to
118 calculate true positive(TP), true negative(TN), false positive(FP), false negative(FN) by
119 comparing with actual output values. After this using the values of TP, TN, FP, FN we need to
120 calculate accuracy, precision, recall, f-measure.

121

122 **5     RESULTS**

123

124 After predicting the values using the model, the performance of model can be evaluated using
125 four metrics namely accuracy, precision, recall, f-measure.

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN} \tag{1}$$

$$Precision = \frac{TP}{TP+FP} \tag{2}$$

$$Recall = \frac{TP}{TP+FN} \tag{3}$$

Where TP = True Positives, TN = True Negatives, FP = False Positives, and FN = False
126 Negatives.
127

```
In [339]: accuracy=num/den*100
          print(accuracy)

          98.24561403508771

In [340]: n=cm[0][0]
          m=cm[0][1]
          recall=n/(n+m)*100
          print(recall)

          100.0

In [341]: p=cm[1][0]
          precision=n/(n+p)*100
          print(precision)

          97.14285714285714

In [342]: fmeasure=(2*recall*precision)/(recall+precision)
          print(fmeasure)

          98.55072463768116
```
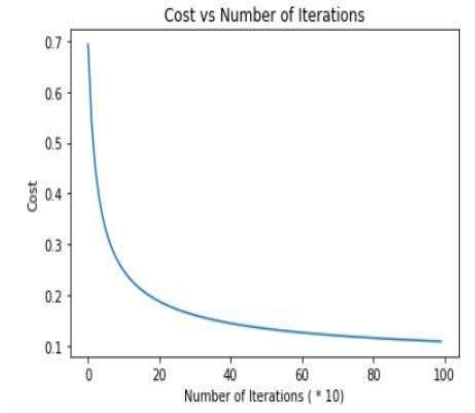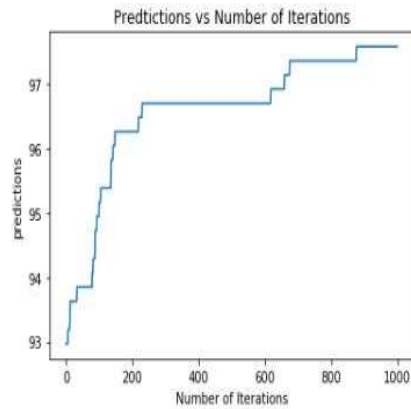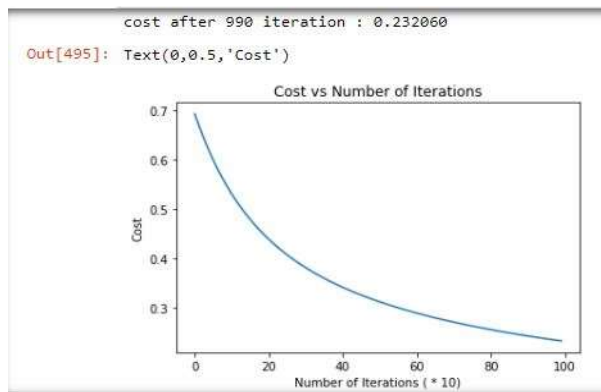
128
129

130 I trained my model with learning rate=0.01 and number of epochs=1000, then I plotted graph
131 between predictions vs number of iterations and number of iterations and cost.

Predictions vs Number of Iterations


Cost vs Number of Iterations

I validated my model using validation data with learning rate=0.001 and number of epochs=1000 and I got cost=0.232060


cost after 990 iteration : 0.232060
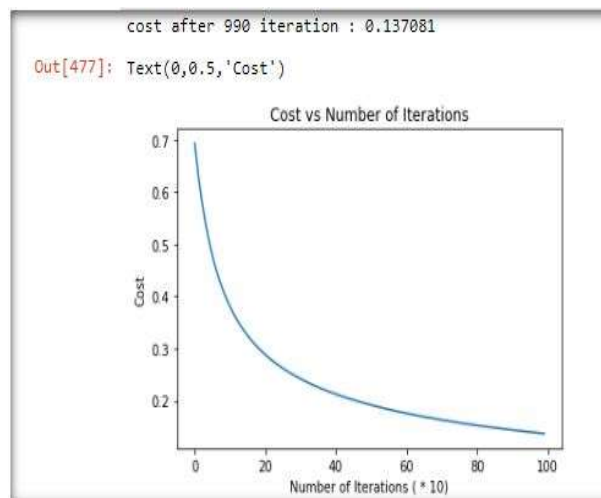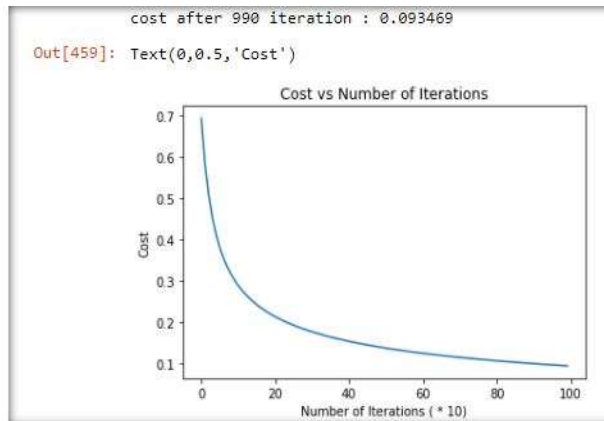Out[495]: Text(0,0.5,'Cost')
Cost vs Number of Iterations

I validated my model using validation data with learning rate=0.003 and number of epochs=1000 and I got cost=0.137081


cost after 990 iteration : 0.137081
Out[477]: Text(0,0.5,'Cost')
Cost vs Number of Iterations

I validated my model using validation data with learning rate=0.006 and number of epochs=1000 and I got cost=0.093469
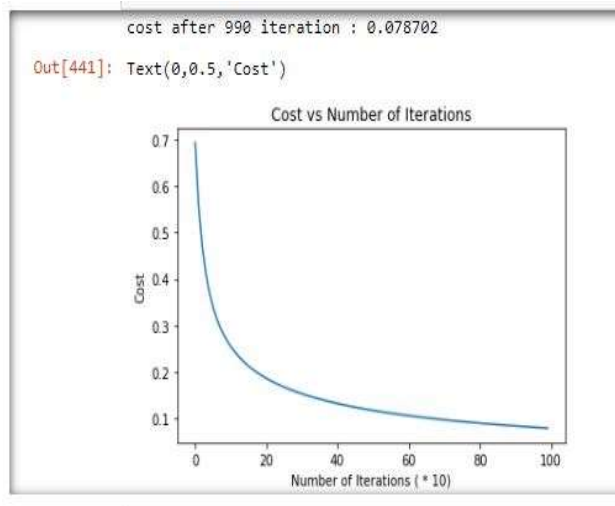
146



147
148

149 I validated my model using validation data with learning rate=0.008 and number of
150 epochs=1000 and I got cost=0.078702
151



152
153
154
155
156 **6       CONCLUSION**
157

158 I successfully trained my model using given cancer dataset and able to validate my dataset as
159 per the graphs shown above. And the following results are drawn.

160

161 Accuracy = 98.24561403508771

162 Precision = 97.14285714285714

163 Recall = 100.0

164 F-measure = 98.55072463768116

165 **References**

166 1) https://machinelearningmastery.com/logistic-regression-tutorial-for-
167 machine-learning/

168    2)  https://en.wikipedia.org/wiki (htt)/Logistic_regression

169    3) https://stats.idre.ucla.edu/r/dae/ordinal-logistic-regression/

170    4) https://ml-cheatsheet.readthedocs.io/en/latest/logistic_regression.html

171    5) https://towardsdatascience.com/logistic-regression-detailed-overview-
172    46c4da4303bc

173    6) https://www.statisticssolutions.com/regression-analysis-ordinal-
174    regression/

175    7) http://ronny.rest/blog/post_2017_08_12_logistic_regression_derivative/