
Predicting Malignant or Benign using Logistic Regression Algorithm

Swati Sajee Kumar

Department of Computer Science Engineering
swatisaj@buffalo.edu #50289560

Abstract

To perform Classification using Logistic Regression on Wisconsin Diagnostic Breast Cancer dataset. The cost minimization is done using gradient descent algorithm. The suspected FNA cells (or predicted output) has been classified into Benign (Class 0) or Malignant (Class 1). The prediction accuracy, precision and recall values were obtained.

1 Introduction

Logistic Regression is an algorithm used for Supervised Learning, which is used for classification purpose. In our project we use Logistic Regression for Binary Classification, that is $y \in \{0,1\}$ where y is the output. For further discussion the symbols and variables used would be as shown in Table 1:

Table 1: Variables used in the report

VARIABLE	FUNCTION
y	Actual Output Expected (0 or 1)
$h_{\theta}(x)$	Predicted output or Hypothesis function (0 or 1) using the training model
θ	Weights
X	Input data
m	Number of samples
n	Number of features
$J(\theta)$	Cost Function, total cost for the m sample dataset

- The hypothesis function h_{θ} lies between 0 and 1, that is $0 \leq h_{\theta} \leq 1$.
- For logistic regression hypothesis function can be written as :
$$h_{\theta}(x) = g(\theta^T x) \dots \dots \dots \text{Eq(1)}$$

where $g(z) = 1/(1+e^{-z})$, on substituting the z value we get
$$g(\theta^T x) = 1/(1+e^{-\theta^T x}) \dots \dots \dots \text{Eq(2)}$$

This is called Sigmoid Function / Logistic Function

- $h_{\theta}(x)$ can be interpreted using two probability functions as shown below:
$$h_{\theta}(x) = P(y=1 | x; \theta) \text{ or } h_{\theta}(x) = P(y=0 | x; \theta)$$

Also, $P(y=1 | x; \theta) = 1 - P(y=0 | x; \theta)$
- Together with the Sigmoid function, probability function, the decision boundary can be created with a threshold = 0.5, that is:
Predict “ $y=1$ ” if $h_{\theta}(x) \geq 0.5$ or
Predict “ $y=0$ ” if $h_{\theta}(x) < 0.5$
- Cost Function / Error Function ($J(\theta)$) for all the samples m , needs to be minimized in

order to predict the hypothesis accurately. Cost function can be represented as follows:

$$J(\theta) = \frac{1}{m} \left(\sum \text{Cost} (h_{\theta}(x), y) \right)$$
 ;where summation runs from $i=1$ to m
 Where $\text{Cost} (h_{\theta}(x), y)$ is cost function for one particular sample data

$$\text{Cost} (h_{\theta}(x), y) = \begin{cases} -\log (h_{\theta}(x)) & \text{if } y=1 \\ -\log (1-h_{\theta}(x)) & \text{if } y=0 \end{cases}$$

 Cost function for one sample data can be written as :

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

 Thus ,

$$J(\theta) = -\frac{1}{m} \left[\sum y \log (h_{\theta}(x)) + (1-y) \log (1-h_{\theta}(x)) \right] \dots\dots\dots \text{Eq(3)}$$

 Furthermore, how we update the weights to minimize the Cost function ($J(\theta)$) and to obtain the hypothesis function $h_{\theta}(x)$ is discussed in Section-4 under Logistic Regression Algorithm.

2 Dataset Description

Wisconsin Diagnostic Breast Cancer(WDBC) dataset was used for training , validation and testing. The dataset contained 569 instances with 32 attributes (ID, diagnosis (B/M), 30 real-valued input features).That is our ‘m’ , number of sample data is 569 and ‘n’ which is number of weights/features is 30.

3 Data Preprocessing

The dataset was divided into three categories : Training dataset (80%) , testing dataset (10 %) and validation dataset (10 %).

3.1 Training Dataset

80 % of the Dataset is randomly chosen to train the model. This data is used to get the weights for minimal cost. Let total data in train set be ‘tr’

3.2 Testing Dataset

10% of the dataset which wasn’t chosen for training is used for testing the algorithm. This will help us to know how accurate our hyperparameters (that is learning rate and weights) are. Let total data in test set be ‘tst’.

3.3 Validation Dataset

10% of the dataset which wasn’t chosen for training or testing the algorithm is used for validation. Validation set is also used like testing data set, but here with this data we change hyperparameters to get accurate results. Let total data in validation set be ‘val’.

Thus, $m = 569 = \text{tr} + \text{tst} + \text{val}$.

4 Logistic Regression Algorithm

Step -1 : Read the data file- using Pandas library

Step -2 : Process the data file – (i)drop the column id/ patient id (ii) Map the labels M/B to 1/0 respectively.

Step -3 : Split the dataframe into train , validation and testing as discussed in Section-3 using sklearn’s `train_test_split`

Step -4 : Normalize the dataset except the first colum. Between the range of 1 to 20. Using the formula –

$$\text{data_norm} = ((\text{data} - \text{data.min}()) / (\text{data.max}() - \text{data.min}())) * 20$$

***** Training step starts from Step-5 to Step-7 *****

Step -5 : Initialize the Learning Rate (lr),also pass the training dataset to the Logistic Regression Model.

100 *Step -6 :Inside the model, initialize the weights and bias, in this project*
101 *weights and bias were initialized to zero.*

102

103 *Step -7 :The loop for gradient descent begins here:*

104

105 **for i in range(10000):**

106

107 $\mathbf{Z} = \boldsymbol{\theta}^T \mathbf{x} + \mathbf{b}$; However we keep the bias as zero always

108

109 • Now we calculate the sigmoid function using the formula below

110

111 **Sigmoid = 1/1+e^{-z}**

112

113 • Now lets calculate the Loss function using the formula obtained from Eq(3), h_{θ} in Eq(3)
114 is the sigmoid here.

115

116 **J($\boldsymbol{\theta}$) = - [y log (sigmoid) + (1-y) log (1- sigmoid)] . mean()**

117

118 • Now we find the partial derivative of the cost function using gradient descent as shown
119 below .

120 • Partial derivative of cost function is written as, $\partial J(\boldsymbol{\theta}) / \partial \theta_j = [\sum_{i=1}^{\text{tr}} (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)}_j]$
121 where i is from 1 to tr (total data in train set) and j is the number iterations used for
122 obtaining the updated weights. *The derivation is shown in Appendix*

123

124 • $\theta_j = \theta_j - \text{lr} [\partial J(\boldsymbol{\theta}) / \partial \theta_j]$ is the formula for updating the weights

125

126 **$\theta_j = \theta_j - \text{lr} [(\text{Sigmoid} - y_{\text{train}}) \cdot x^T]$**

127

128 • In the above equation we are updating the weights using the formula $\theta_j = \theta_j - \text{lr} [\sum_{i=1}^{\text{tr}} (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)}_j]$ where h_{θ} is our sigmoid function , and y is the training data set.
129 Also 'lr' here is the learning rate.

130

131 • This loop runs until i – 10000

132

133 ***** Testing Model in Step-8 *****

134

135 *Step -8 :Create another module for prediction, where we would pass the*
136 *training / validation dataset and the updated weights obtained from training*
137 *the module from Step 7*

138

139 • Use the genesis equation and find the Z , the weights here are the ones which we obtained
140 from training the module

141

142 **$\mathbf{Z} = \boldsymbol{\theta}^T \mathbf{x} + \mathbf{b}$**

143

143 • Now find the sigmoid function using the updated weights

144

145 **Sigmoid = 1/1+e^{-z}**

146

146 • Now start the loop for prediction

147

148 **for i in range(size of sigmoid)**

149 **if(sigmoid <0.5):**

150 **y_predict = 0**

151 **else:**

152 **y_predict = 1**

153

153 • The above step classifies the output value to Benign (y_predict = 0)
154 or Malignant (y_predict = 1) on the basis of sigmoid value 0.5

155 The True Positive, True Negative , False Positive and False Negative values
 156 are checked. The results for it are mentioned in Section-5 , Results and
 157 Discussion section.

158

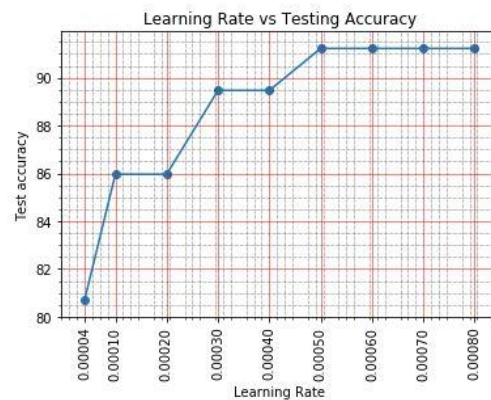
159 5 Results and Discussion

- 160 • The number of iterations for every epoch was 10000
- 161 • The length of training dataset was 455
- 162 • The length of testing dataset was 57
- 163 • The length of validation dataset was 57
- 164 • The different values of Learning Rate used to obtain the maximum
 165 testing accuracy and minimum loss are as follows:

166 lr_list=[0.00004,0.0001,0.0002,0.0003,0.0004,0.0005,0.0006,0.0007,0.0008]
 167

168 At the end of Results and Discussion we have mentioned for which value of
 169 LR we obtained the best Testing Accuracy and Loss Function.

170 5.1 Learning Rate vs Testing Accuracy

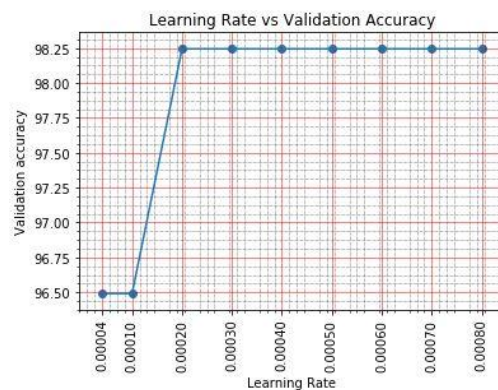


171

172 Figure 1: Learning Rate vs Testing Accuracy

173 From Fig 1. , we can see that as the learning rate increases from 0.00004 to 0.0005 the
 174 Testing Accuracy increases. Also we can notice that after 0.0005 the testing accuracy
 175 becomes constant. Furthermore, if we keep changing the values of LR it is more likely that
 176 after a point the testing accuracy will decrease again. This is because of the convex nature of
 177 the graph, ones we reach the local minima the accuracy is the highest, and it again decreases
 178 after a certain point.

179 5.2 Learning Rate vs Validation Accuracy



180

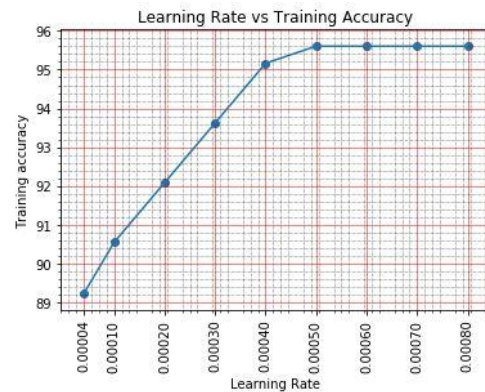
181 Figure 2: Learning Rate vs Validation Accuracy

182

183 We could see that the validation accuracy increases linearly from 0.0001 to 0.0002, after
184 which the accuracy becomes constant.

185 5.3 Learning Rate vs Training Accuracy

186

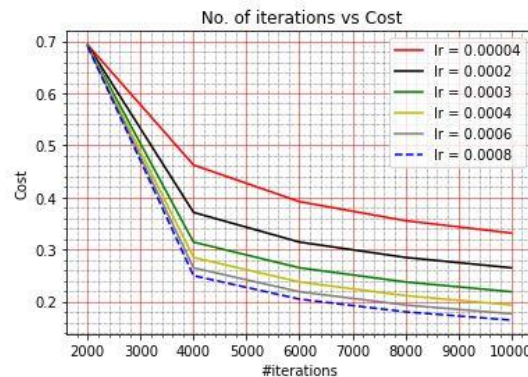


187

188 Figure 3: Learning Rate vs Training Accuracy

189 We can see that the training accuracy increases as we go from Learning Rate = 0.00004 to
190 0.0005 almost linearly, later as seen in testing accuracy's graph here as well the accuracy
191 becomes constant after a 0.0005. . Furthermore, if we keep changing the values of LR it is
192 more likely that after a point the training accuracy will decrease again. This is because of the
193 convex nature of the graph, ones we reach the local minima the accuracy is the highest, and
194 it again decreases after a certain point.

195 5.4 Cost vs Number of Iterations



196

197 Figure 4: Cost vs Number of Iterations

198

199 Here we check the Cost/ Error function with respect to the number of iterations. We can note
200 that as the number of iterations increases from 2000 to 10000 keeping LR as constant the
201 Cost decreases. Also we can note that for LR=0.0008 the cost decreases rapidly. Also we can
202 see that we get a minimum Loss/Cost function at LR = 0.0008 and iterations = 10000

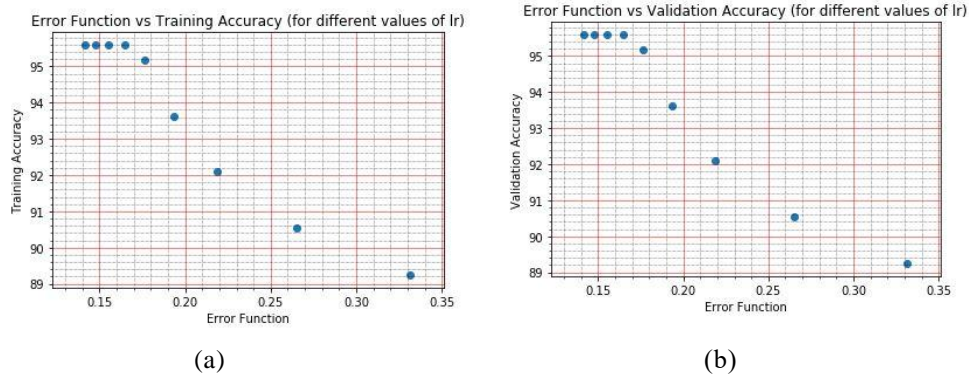
203

204

205

206

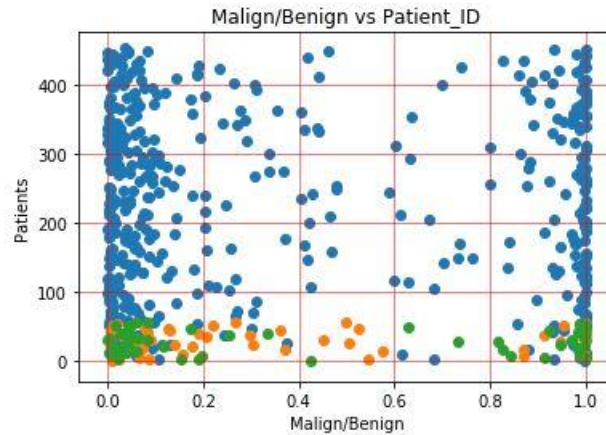
207 5.5 Error function vs Accuracy



210 Figure 5: (a) Error function vs Testing Accuracy (b) Error function vs Validation Accuracy

212 The above two graph depicts that the Accuracy in general decreases if the Error
 213 Function/Cost/Loss function increases. Thus our aim needs to include not only an accurate
 214 value , but accuracy with minimum Error. Both of them are inversely proportional to each
 215 other.

216 5.6 Predicting the output for testing data



218 Figure 4: Predicting the output for training data

219 There were 57 data points in testing data, this scattered image was obtained before we
 220 classified our sigma values as malignant or benign. We can notice that the sigmoid function
 221 in prediction model along with the weights calculated in training section got us allocate the
 222 features between 0 and 1 , as shown above. It is after this that we classify the points lying
 223 above 0.5 as malignant and below 0.5 as benign.

224 **5.7 Accuracy, Precision and Recall** were calculated for Training dataset. Loss Function was
 225 also noted down. Along with training and validation accuracy for different values of LR as
 226 shown in Table 2.

$$227 \text{ Accuracy} = (TP + TN) / (TP + TN + FP + FN)$$

$$228 \text{ Precision} = TP / (TP + FP)$$

$$229 \text{ Recall} = TP / (TP + FN)$$

230 Where TP = True Positive , TN = True Negative, FP = False Positive and FN = False
 231 Negative

232 The best accuracy and minimum loss function was obtained at **LR = 0.0008** and **iterations =**
 233 **10000** the corresponding values are mentioned in the tabular column below:

234 Table 2: Precision Parameters for LR = 0.0008

LR	TP	TN	FP	FN	Testing Accuracy	Testing Precision	Testing Recall	Training Accuracy	Validation Accuracy	Loss Function
0.00004	12	34	7	4	80.701%	63.1578%	75.0%	89.23%	96.49%	33.1613%
0.0001	13	36	5	3	85.964%	72.221%	81.25%	90.549%	96.491%	26.493%
0.0002	13	36	5	3	85.964%	72.225%	81.25%	92.087%	98.24%	21.8956%
0.0003	13	38	3	3	89.47%	81.25%	81.25%	93.626%	98.245%	19.3533%
0.0004	13	38	3	3	89.47%	81.25%	81.25%	95.16%	98.24%	17.66%
0.0005	14	38	3	2	91.22%	82.35%	87.5%	95.60%	98.24%	16.45%
0.0006	14	38	3	2	91.22%	82.35%	87.5%	95.60%	98.24%	15.2278%
0.007	14	38	3	2	91.22%	82.35%	87.5%	95.60%	98.24%	14.7844%
0.0008	14	38	3	2	91.228 %	82.3529%	87.5%	95.60439%	98.2456%	14.179%

235 7 Conclusion

237 The best epoch for which the Testing Accuracy was maximum and the loss function was
 238 minimum was for **learning rate = 0.0008** and **iterations = 10000**. The Test Accuracy
 239 obtained was **91.228 %** , Precision obtained was **82.3529 %** and Recall value was **87.5%**.

240 Acknowledgments

241 I am extremely grateful to Professor Sargur Srihari for elucidating all the necessary concepts
 242 related to Logistic Regression. I am also thankful to all the Teaching Assistants, their recitation
 243 sessions have been really helpful.

244 References

- 245 [1] Andrew Ng's course on Youtube :
 246 [https://www.youtube.com/watch?v=TTdcc21Ko9A&list=PLLssT5z_DsK-](https://www.youtube.com/watch?v=TTdcc21Ko9A&list=PLLssT5z_DsK-h9vYZkQkYNWcItqhlRJLN&index=36)
 247 [h9vYZkQkYNWcItqhlRJLN&index=36](https://www.youtube.com/watch?v=TTdcc21Ko9A&list=PLLssT5z_DsK-h9vYZkQkYNWcItqhlRJLN&index=36)
 248 [2] CSE474LECB: Machine Learning – Week 3: Linear Regression --→ Gradient Descent
 249 [3] Regarding TP, TN, FP and FN I referred [https://stackoverflow.com/questions/31324218/scikit-](https://stackoverflow.com/questions/31324218/scikit-learn-how-to-obtain-true-positive-true-negative-false-positive-and-fal)
 250 [learn-how-to-obtain-true-positive-true-negative-false-positive-and-fal](https://stackoverflow.com/questions/31324218/scikit-learn-how-to-obtain-true-positive-true-negative-false-positive-and-fal)
 251 [4] PRML, Section 4.3.2 for Logistic Regression

252
 253
 254
 255
 256
 257
 258
 259
 260
 261
 262

APPENDIX

Partial derivation of Cost Function w.r.t Weights :

AUTHOR : SWATI SAJEE KUMAR

We know,

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m y^i \log(h_{\theta}(x^i)) + (1-y^i) \log(1-h_{\theta}(x^i)) \quad (1)$$

where $J(\theta) \rightarrow$ cost

$m \rightarrow$ no. of samples

$x \rightarrow$ i/p data

$h_{\theta}(x^i) \rightarrow$ sigmoid function

$y \rightarrow$ Actual o/p

Also we know that,

$$h_{\theta}(x^i) = \frac{1}{1+e^{-\theta^T x^i}} \Rightarrow \log(h_{\theta}(x^i)) = \log \frac{1}{1+e^{-\theta^T x^i}}$$

$$\Rightarrow \log(h_{\theta}(x^i)) = -\log(1+e^{-\theta^T x^i}) \quad (2)$$

Substituting this in equation (1) we get

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \left[-y^i \log(1+e^{-\theta^T x^i}) + (1-y^i) (-\theta^T x^i - \log(1+e^{-\theta^T x^i})) \right]$$

which can be simplified to,

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \left[y_i \theta^T x^i - \theta^T x^i - \log(1+e^{-\theta^T x^i}) \right] \quad (3)$$

$$-\theta^T x^i - \log(1+e^{-\theta^T x^i}) = -\int \log e^{\theta^T x^i} e^{-\theta^T x^i} d(\theta^T x^i)$$

$$= -\log(1+e^{\theta^T x^i}) + \log(1+e^{-\theta^T x^i})$$

Now if we substitute this value in equation ③ we get

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \left[\underbrace{y_i \theta^T x^i}_{\textcircled{A}} - \log \underbrace{\left(1 + e^{\theta^T x^i} \right)}_{\textcircled{B}} \right]$$

Now let's compute partial derivative w.r.t θ_j for element A;

$$\frac{\partial}{\partial \theta_j} y_i \theta^T x^i = y_i x_j^i$$

for element B;

$$\begin{aligned} \frac{\partial}{\partial \theta_j} \log(1 + e^{\theta^T x^i}) &= \frac{x_j^i e^{\theta^T x^i}}{1 + e^{\theta^T x^i}} \\ &= x_j^i h_0(x^i) \end{aligned}$$

$$\therefore \frac{\partial J}{\partial \theta_j} = -\frac{1}{m} \sum_{i=1}^m (h_0(x^i) - y_i) x_j^i$$

Hence Derived