

Assignment – 2

Swapnil Siddharth Course: 1CSD1 Student ID: 21230014

Algorithm: Logistic Regression

Logistic Regression is a classification technique which is used to predict distinct set of classes based on given observations. As linear regression technique produces continuous output values, logistic regression transforms those output values into a probability value using a sigmoid function. These probabilistic values can then further be represented into two or more distinct classes.

Most machine learning algorithms pertain to certain assumptions, Naive Bayes makes assumptions of conditional independence of features, K-NN makes assumptions that the query point resembles the properties neighborhood points. The main assumption we make in the logistic regression algorithm is that our classes are almost linearly separable or in some cases perfectly separable. This means that if there are two class points, we can draw a hyperplane which can separate the positive class from negative class.

Let us suppose that π is the decision surface plane which best separates the two class points, X is the feature vector, W is normal (which can be also called as weights) to the plane and B is the intercept. Then, $W^T X + B = 0$ is the equation of the plane, where W and X are vectors and B is a scalar.

Then we can pass this equation to the sigmoid function, which intakes any real valued number and provides an output between the range of 0 and 1. Hence on applying the sigmoid function, we get: $Y = \sigma(W^T X + B)$, where Y is the prediction or the probability of event occurrence. If the probability is greater than 0.5, we classify it as 1, if probability is less than 0.5, we classify it as 0.

Our task in model building of Logistic Regression is to find the best W and B which corresponds to the plane such that can best separate positive points from negative points, and hence it can help us in predicting the class labels more accurately.

Design Schemes:

Linear Model

Logistic Regression is based on the fact that the target variable is categorical in nature and has only two possible classes. And using this analogy, it can apply probabilistic approach to find the occurrence of any event. We initially calculate the linear combination of dependent and independent variables, which is basically the linear regression model, which will provide us continuous output values. On that model we apply the sigmoid function which will give us a probability.

The decision surface in Logistic Regression is a hyperplane or line. And hence any point which is on the direction of the normal to the plane will be positive and the points which is opposite direction of the normal will be negative. Finding the right plane is the main task in building model such that it maximizes sum of correctly classified points(positive) and misclassified points(negative). But here, outliers can highly impact the plane, and consequently the sum of sum of correctly classified points(positive) and misclassified points(negative). Thus, we can do squashing using sigmoid function, where the maximum value is 1 and minimum is 0.

Sigmoid Function

The Sigmoid function is creating an S-shaped curve and is often called as squashing function as it can intake any real valued number and its output range is between 0 and 1. The advantage of this function is that it can input a very large positive number or a very large negative number, but the function will give output in between 0 and 1. If the curve tends towards positive infinity, the predicted value will be close to 1, consequently if the curve tends towards negative infinity, the predicted value will be almost 0.

Since outliers can highly impact the decision surface(plane) and our summation of correctly classified points(positive) and misclassified points(negative) in our linear model, performing squashing using sigmoid function can solve this issue, where the maximum value is 1 and minimum is 0. If the value is small, the curve will grow almost linearly and beyond 1 it will taper off and similarly beyond 0, it will taper off the other side.

The advantage of sigmoid function is that it has a probabilistic interpretation and can be differentiable. From probabilistic point of view, if a point is on very far the direction normal to the plane from the plane, the probability that it's a positive point reaches close to 1 and if very far from opposite to the direction of plane it will be close to 0.

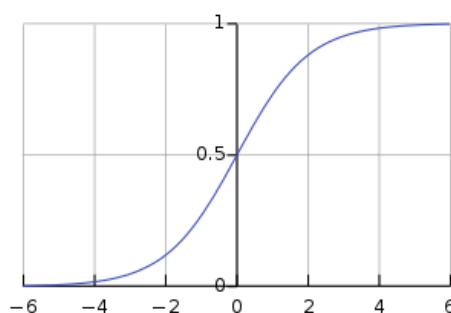


Figure 1: Sigmoid Function logistic curve

Source: Wikipedia

Cost Function

Cost function is basically to represent the errors present in our model building. It can help us to calculate how correctly our model is predicting the labels as compared to the actual class labels in the train dataset. If the cost function is more, the accuracy of predicting the actual class will reduce, consequently, if the cost function is less, the accuracy of predicting the actual class will increase. In Logistic Regression, we use cross-entropy cost function which is

called log loss. The advantage of using cross-entropy is that it uses the concept of monotonic functions to create the cost functions.

A function $f(x)$ is said to be monotonic function such that if x increases $f(x)$ also increases and vice versa. Thus, this helps us in creating smooth cost function graphs which helps us in calculating the gradient and reducing the cost.

Thus, the cost function for logistic regression can be represented as:

$$cost = \frac{1}{n} \sum_{i=1}^n Y_i \cdot \log(p(y^i)) + (1 - y^i) \cdot \log(1 - p(y^i))$$

Here, y^i is the actual class label and $\log(p(y^i))$ is the probability of that class label.

Gradient Descent

Gradient descent is an optimization algorithm which helps us in finding the global or local minima of the function. It is inherently done to reduce the cost function. This optimization technique helps the model to learn the gradient in other words the direction which model should take to reduce errors. On plotting the graph of cost function, the main task of gradient descent is to calculate the weight and bias such that, as it starts from the top of curve with every step iteration it moves down in the direction which minimizes the cost function straight away. This is done iteratively such that we can reach the bottom of the graph which becomes our local minima. Here, the learning rate decides the step size the gradient should make in each iteration to reach the minima. The learning rate, which is low, where we are recalculating the minima very frequently, is much preferred as it can more precisely move in the direction of negative gradient. The whole process is integral, and the advantage of sigmoid function used in Logistic Regression is that its derivative is easy to calculate. Hence, we can build an equation for cost function derivative as:

$$W = x(p - a)$$

where, W = derivative of cost with respect to weights

p = prediction by the model

a = actual class label (0 or 1)

x = feature vector

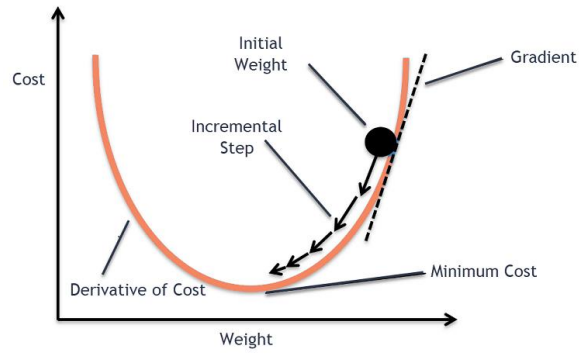


Figure 2: Gradient Descent

Source: Clairvoyant

Tests and Results:

Using the above-mentioned designed scheme, an implementation has been done to build the logistic regression classifier. The input file(wildfires) is a text file, upon which data cleaning and pre-processing has been performed such that the unneeded spaces have been removed and proper header to each feature column has been provided. Then using the Label Encoder provided with preprocessing under scikit-learn package has been used to turn the categorical value of yes and no for fire into 1 and 0. The training data has been randomly shuffled and divided dataset such as two-third of dataset is used for training and one-third is used for testing. The training data is then fit into the logistic regression model created as well as the built-in logistic regression classifier provided by scikit-learn package. The predictions of class labels are done for each model iteratively ten times and the accuracy score are calculated in each iteration. Afterwards, a mean accuracy score is calculated for implemented model as well as the reference model. Upon training the implementation of the logistic regression model on the train dataset and using the model to predict the class labels, we get an accuracy score of 72.35% as compared to 86.61% of the reference logistic regression classifier of scikit-learn.

```

Iteration Number: 1
Implemented Logistic Regression Classifier model accuracy: 69.11764705882352
Scikit-Learn Logistic Regression Classifier model accuracy: 88.23529411764706

Iteration Number: 2
Implemented Logistic Regression Classifier model accuracy: 86.76470588235294
Scikit-Learn Logistic Regression Classifier model accuracy: 88.23529411764706

Iteration Number: 3
Implemented Logistic Regression Classifier model accuracy: 63.23529411764706
Scikit-Learn Logistic Regression Classifier model accuracy: 83.82352941176471

Iteration Number: 4
Implemented Logistic Regression Classifier model accuracy: 67.64705882352942
Scikit-Learn Logistic Regression Classifier model accuracy: 85.29411764705883

Iteration Number: 5
Implemented Logistic Regression Classifier model accuracy: 64.70588235294117
Scikit-Learn Logistic Regression Classifier model accuracy: 86.76470588235294

Iteration Number: 6
Implemented Logistic Regression Classifier model accuracy: 75.0
Scikit-Learn Logistic Regression Classifier model accuracy: 85.29411764705883

Iteration Number: 7
Implemented Logistic Regression Classifier model accuracy: 79.41176470588235
Scikit-Learn Logistic Regression Classifier model accuracy: 91.17647058823529

Iteration Number: 8
Implemented Logistic Regression Classifier model accuracy: 63.23529411764706
Scikit-Learn Logistic Regression Classifier model accuracy: 89.70588235294117

Iteration Number: 9
Implemented Logistic Regression Classifier model accuracy: 72.05882352941177
Scikit-Learn Logistic Regression Classifier model accuracy: 85.29411764705883

Iteration Number: 10
Implemented Logistic Regression Classifier model accuracy: 82.35294117647058
Scikit-Learn Logistic Regression Classifier model accuracy: 82.35294117647058

Mean Accuracy of Implemented Logistic Regression Classifier model accuracy: 72.35294117647058
Mean Accuracy of Scikit-Learn Logistic Regression Classifier model accuracy: 86.61764705882352

```

Figure 3: Accuracy Results

Conclusions and observations:

After training the implemented logistic regression model on the train dataset and using the model to predict the class labels, we get an accuracy score of 72.35% as compared to 86.61% of the reference logistic regression classifier of scikit-learn. During initially training and testing of the model, we observed that the learning rate if increased and the number of iterations decreased, we are getting a lower accuracy score. Thus, upon experimenting with the learning rate and number of iterations, we find the ideal values of learning rate and number of iterations to be 0.0001 and 10000 respectively. This gives us a fairly good accuracy performance which is slightly less than the accuracy of the model trained with scikit-learn logistic regression classifier. An output file has also been created which stores the output predicted and actual values in a csv file.

References:

- How Does the Gradient Descent Algorithm Work in Machine Learning?
[<https://www.analyticsvidhya.com/blog/2020/10/how-does-the-gradient-descent-algorithm-work-in-machine-learning/>]
- Sigmoid function [https://en.wikipedia.org/wiki/Sigmoid_function]
- Scikit-Learn GitHub [https://github.com/scikit-learn/scikit-learn/blob/main/sklearn/linear_model/_logistic.py]
- Logistic Regression [https://en.wikipedia.org/wiki/Logistic_regression]
- Cross entropy [https://en.wikipedia.org/wiki/Cross_entropy]
- Understanding Logistic Regression in Python
[<https://www.datacamp.com/community/tutorials/understanding-logistic-regression-python>]

Appendix:

```

1  #Importing required libraries for the assignment code
2
3  import pandas as pd #for data-preprocessing
4  import numpy as np #for creating nd-arrays and matrix multiplications
5  from sklearn.model_selection import train_test_split #for splitting dataset into
   test and train data
6  from sklearn.linear_model import LogisticRegression #for using Logistic Regression
   classification and compare it's accuracy with bulid classifier
7  from sklearn import metrics #for accuracy calculations of
8  from sklearn import preprocessing #for importing LabelEncoder
9
10 wildfires_data = open('wildfires.txt', 'r')
11
12
13 #Creating a function which removes any unneeded spaces in the text dataset
14 def remove_spaces(input_file):
15     elements = []
16     for line in input_file:
17         element = []
18         element = line.split()
19         elements.append(element)
20     return elements
21
22 #Creating the required column names of the dataset
23 column_names = ['fire', 'year', 'temp', 'humidity', 'rainfall', 'drought_code',
   'buildup_index', 'day', 'month', 'wind_speed']
24
25 #Applying the function to remove spaces and creating the dataframe using pandas and
   giving the appropriate column names
26 wildfires_df = pd.DataFrame(remove_spaces(wildfires_data), columns = column_names)
27
28 #Dropping the first row of the dataframe as it contains column names
29 wildfires_df = wildfires_df.iloc[1:]
30
31 #Creating feature columns
32 features_columns = ['year', 'temp', 'humidity', 'rainfall', 'drought_code',
   'buildup_index', 'day', 'month', 'wind_speed']
33
34 #creating labelEncoder
35 label_encoder = preprocessing.LabelEncoder()
36
37 #Label encoding the fire column
38 wildfires_df['fire'] = label_encoder.fit_transform(wildfires_df['fire'])
39 #print(wildfires_df)
40
41 #Creating X and y dataframes with required features and classes from wildfires
   dataset
42 X = wildfires_df[features_columns]
43 y = wildfires_df['fire']
44
45 #Creating Logistic Regression Classifier
46 class LogisticRegressionClassifier:
47
48     #Creating init method and initializing the values 0.001 and 1000 to learning
   rate and number of iterations for gradient descent respectively
49     def __init__(self, learning_rate=0.001, number_of_iterations=1000):
50         self.learning_rate = learning_rate
51         self.number_of_iterations = number_of_iterations
52         #Creating weights and bias as None initially, which later need to be
   calculated
53         self.weights = None
54         self.bias = None
55
56     #Creating the sigmoid function
57     def sigmoid(self, x):
58         return 1 / (1 + np.exp(-x))
59
60     #Creating method to calculate gradient descent and updating weights and bias
61     def gradient_descent(self, n_sample, p_score, actual, Z):
62
63         #Calculating the derivative of weights and bias
64         dw = (1 / n_sample) * np.dot(Z.T, (p_score - actual))

```

```

65         db = (1 / n_sample) * np.sum(p_score - actual)
66
67         #Adjusting the weights and bias with respect to learning rate
68         self.weights = self.weights - self.learning_rate * dw
69         self.bias = self.bias - self.learning_rate * db
70
71     #Creating the fit method which will take training dataset and labels as parameters
72     def fit(self, train_data, target_labels):
73
74         #Initializing the parameters
75         number_of_samples = train_data.shape[0]
76         number_of_features = train_data.shape[1]
77         self.weights = np.zeros(number_of_features)
78         self.bias = 0
79
80         #Running iterations for building linear model to fit into sigmoid function
81         #and updating weights and bias using gradient descent function
82         for _ in range(self.number_of_iterations):
83
84             #Creating a linear model with y and dot product of weights and training
85             #features data adding bias
86             linear_model = np.dot(train_data, self.weights) + self.bias
87
88             #Applying the build sigmoid function to the linear model to find the
89             #probabilistic scores (approximations) of y
90             probability_scores = self.sigmoid(linear_model)
91
92             #Calling Gradient Descent function
93             self.gradient_descent(number_of_samples, probability_scores,
94                                   target_labels, train_data)
95
96     def predict_class(self, X):
97
98         #Creating a linear model with y and dot product of weights and x adding bias
99         linear_model = np.dot(X, self.weights) + self.bias
100
101         #Applying the build sigmoid function to the linear model to find the
102         #approximation of y
103         probability_scores = self.sigmoid(linear_model)
104
105         #Using list comprehension to predict classes (values > 0.5 as 1 and values <
106         #0.5 as 0)
107         predicted_classes = [1 if i > 0.5 else 0 for i in probability_scores]
108         return np.array(predicted_classes)
109
110
111     #Defining the accuracy function to calculate the accuracy of the build model
112     def accuracy_score(actual_class, predicted_class):
113         accuracy_score = np.sum(actual_class == predicted_class) / len(actual_class)
114         return accuracy_score
115
116     #Building list to store accuracy for 10 iterations and calculating mean accuracy score
117     build_model_accuracy_list = []
118     scikitlearn_model_accuracy_list = []
119
120     for i in range(10):
121
122         print("Iteration Number:", i+1)
123
124         #Splitting the dataset into train and test datasets
125         X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33,
126                                                             shuffle = True)
127
128         #Converting the train and test datasets into numpy nd-array
129         X_train, X_test, y_train, y_test = X_train.values , X_test.values ,
130                                             y_train.values , y_test.values
131
132         #Converting the nd-array of test and train dataset into float values in order to
133         #facilitate matrix multiplications
134         X_train, X_test, y_train, y_test = X_train.astype(float) , X_test.astype(float)

```



```

128     , y_train.astype(float) , y_test.astype(float)
129
130 #Creating a classifier object using the implemented(built) Logistic Regression
131 classifier and setting parameters 0.0001 and 10000 for learning rate and number
132 of iterations
133 logistic_regression_build_model =
134 LogisticRegressionClassifier(learning_rate=0.0001, number_of_iterations=10000)
135
136 #Calling the fit function to train the model with the train datasets
137 logistic_regression_build_model.fit(X_train, y_train)
138
139 #Predicting the class on test dataset using the model
140 build_predictions = logistic_regression_build_model.predict_class(X_test)
141
142 #Calculating the accuracy score with accuracy function
143 build_model_accuracy = (accuracy_score(y_test, build_predictions))*100
144
145 #Printing the accuracy in each iteration
146 print("Implemented Logistic Regression Classifier model accuracy:",
147       build_model_accuracy)
148
149 #Appending the accuracy of each iteration in the respective list
150 build_model_accuracy_list.append(build_model_accuracy)
151
152 #Creating Logistic Regression classifier object using Sk-learn Logistic
153 Regression classifier
154 sklearn_logistic_regression_model = LogisticRegression(solver='lbfgs',
155 max_iter=1000)
156
157 #Training Logistic Regression Classifier
158 logisticregression_clf = sklearn_logistic_regression_model.fit(X_train,
159 y_train)
160
161 #Predicting the response from test dataset
162 sklearn_model_predictions = logisticregression_clf.predict(X_test)
163
164 sklearn_model_accuracy = (metrics.accuracy_score(y_test,
165 sklearn_model_predictions))*100
166 print("Scikit-Learn Logistic Regression Classifier model accuracy:",
167       sklearn_model_accuracy, "\n")
168 sklearn_model_accuracy_list.append(sklearn_model_accuracy)
169
170 print("Mean Accuracy of Implemented Logistic Regression Classifier model accuracy:
171 ", np.mean(build_model_accuracy_list))
172
173 print("Mean Accuracy of Scikit-Learn Logistic Regression Classifier model accuracy:
174 ", np.mean(sklearn_model_accuracy_list))
175
176 #Creating dictionary of actual and predicted labels of model and storing it in an
177 output file
178 actual_labels_and_predicted_labels = pd.DataFrame({'Actual Label': y_test,
179 'Predicted Label': build_predictions})
180 actual_labels_and_predicted_labels.to_csv("label_outputs.csv", index = False)

```